

计算数学丛书

# 对称矩阵计算

蒋尔雄 编著

上海科学技术出版社

51.8  
840

计算数学丛书

# 对 称 矩 阵 计 算

蒋 尔 雄 编 著



上海科学技术出版社

8510301

Dt06/12

计算数学丛书

对称矩阵计算

蒋尔雄 编著

上海科学技术出版社出版

(上海瑞金二路 450 号)

新华书店上海发行所发行 上海商务印刷厂印刷

开本 787×1092 1/32 印张 5.375 字数 115,000

1984 年 11 月第 1 版 1984 年 11 月第 1 次印刷

印数: 1-11,500

统一书号: 13119·1184 定价: 0.75 元

1080109

## 出版说明

《计算数学丛书》是为了适应计算数学和计算机科学的发展，配合高等院校计算数学教学的需要而组织的一套参考读物。读者对象主要是高等院校数学系和计算机科学系的学生、研究生，亦可供高等院校数学系和计算机科学系的教师以及工矿企业、科研单位从事计算工作的技术人员参考。

本丛书向读者介绍近代计算方法的一些主要进展及其适用范围和实用效果。每种书集中介绍一个专题，针对本专题的近代发展作综合性的介绍，内容简明扼要，重点突出，有分析，有评价，力图使读者对该专题的动向和发展趋势得到一个完整的了解。

本丛书已拟定的选题计有：《线性代数与多项式的快速算法》、《数论变换》、《数值有理逼近》、《矩阵特征值问题》、《索伯列夫空间引论》、《计算组合数学》、《样条与插值》、《有限条形法》、《广义逆矩阵及其计算方法》、《非线性方程迭代解法》、《奇异摄动中的边界层校正法》、《沃尔什函数理论与应用》、《多项式最佳逼近的实现》、《坏条件常微分方程数值解》、《误差分析》、《最小二乘问题的数值解法》、《板壳问题非协调方法》、《外推法及其应用》、《Monte Carlo 方法》、《差分格式理论》、《高维偏微分方程数值解》等二十余种，于一九八〇年初起陆续出版。

《计算数学丛书》编辑委员会

主 编

李 荣 华

编 委

冯果忱 李岳生 李荣华 吴文达 何旭初

苏煜城 胡祖炽 曹维潞 雷晋干 蒋尔雄

矩阵计算是近代计算方法中发展最早、最快的领域。尤其是对称矩阵计算是最成熟的部分。

本书比较深入地介绍了解线性代数方程组的共轭斜量法和 Lanczos 方法, 解矩阵特征值问题的 QL 方法和 Lanczos 方法, 以及与这些方法密切相关的对称三对角矩阵的理论。这些方法是经过实践考验的, 认为是有效的方法。其中解线性代数方程组的共轭斜量法和解特征值问题的 Lanczos 方法, 都是五十年代初期就提出来了, 经历了二十年曲折的历史, 获得了理论上和实践上的支持, 直到七十年代才站住了脚。

读者通过本书, 一方面可以更好地掌握这些方法, 应用于实际, 也可以了解到目前在对称矩阵计算理论上达到的深度; 并且会发现, 尽管对称矩阵计算这个领域, 在计算方法各个领域中, 相对来说是比较成熟的, 但仍然有不少问题有待解决。

本书的初稿曾是复旦大学计算数学专业 77 届、78 届学生选修课的讲稿, 在教学过程中, 发现学生对这些内容还是比较感兴趣的。现在写成此书, 作为大学计算数学专业大学生的选修课和研究生课程的教材也是适当的。

本书中有很多素材来自美国 California 大学 Berkeley 分校数学系和计算机科学系教授 B. N. Parlett 的书 “The Symmetric Eigenvalue Problem”, 记得在 1979 年那书正式出版之前, Parlett 教授把这书的打印稿给我阅读, 使我能及时掌握有关的理论。在此向 Parlett 教授表示衷心的感谢。

本书的其它材料，来自书后所列的参考文献和自己的心得。由于水平有限，错误和不妥之处在所难免，望读者批评指正。

蒋尔雄

于复旦大学 1982年10月

## 记号说明

### 本书规定

1. 用大写英文字母表示矩阵, 用黑体小写英文字母表示向量, 用小写英文字母或希腊字母表示标量.

2.  $A^T$ ,  $x^T$  表示  $A$ ,  $x$  的转置;  $A^*$ ,  $x^*$  表示  $A$ ,  $x$  的转置共轭.

3. 列向量  $x = (x_1, x_2, \dots, x_n)^T$ ,  $y = (y_1, y_2, \dots, y_n)^T$  的内积

$$(x, y) = \sum_{i=1}^n \bar{x}_i y_i = x^* y.$$

4. 向量  $x$  的范数  $\|x\|$  表示  $x$  的  $l_2$  范数, 矩阵  $A$  的范数

$$\|A\| = \sup_{x \neq 0} \|Ax\| / \|x\|.$$

5.  $e_i$  表示单位矩阵

$$\begin{pmatrix} 1 & & 0 \\ & 1 & \\ 0 & & \ddots \\ & & & 1 \\ & & 0 & & \end{pmatrix}$$

的第  $i$  列.

6.  $q_1, q_2, \dots, q_j$  是  $j$  个  $n$  维列向量, 记  $[q_1, q_2, \dots, q_j]$  表示由这  $j$  个向量, 按照所给的次序组成的矩阵.

$\{q_1, q_2, \dots, q_j\}$  表示由这  $j$  个向量所成的线性空间.

7. 没有特殊声明, 矩阵、向量、常数全是指实矩阵、实向量、实数; 任意矩阵, 意指任意实矩阵.



# 目 录

引言

记号说明

<b>第1章 共轭斜量法</b>	1
§1 斜量法	1
§2 多步斜量法	10
§3 共轭斜量法	19
§4 不完全分解、预处理共轭斜量法	27
<b>第2章 对称三对角矩阵</b>	33
§1 Jacobi 矩阵	33
§2 对称三对角矩阵的唯一归化定理	36
§3 对称三对角矩阵的极值性质	48
§4 Thompson-McEntegert-Paige 公式和特征值反问题	52
§5 解对称线性代数方程组的 Lanczos 算法	63
<b>第3章 解特征值问题的 QL 方法</b>	76
§1 QL 方法的一般性质	76
§2 用于对称三对角矩阵时的 QL 方法的性质	83
§3 带 Rayleigh 商位移的 QL 方法	97
§4 带 Wilkinson 位移的 QL 方法	103
<b>第4章 解特征值问题的 Lanczos 算法</b>	111
§1 近似不变子空间	112
§2 Lanczos 算法	124
§3 Kaniel-Paige-Saad 理论	131
§4 在有限精度运算下的 Lanczos 算法	143
<b>参考文献</b>	158

共轭斜量法

共轭斜量法 (conjugate gradient method) 是解系数矩阵为对称正定的线性代数方程组的一种方法, 它的产生在五十年代初期, 参见参考文献[1]; 经过几十年的考验, 现在公认为它是一种好方法, 可详阅[2]、[3]、[4]. 为了深入了解共轭斜量法, 我们从最佳逼近的观点来介绍共轭斜量法. 另外我们也介绍一下, 近年来产生的, 引起广泛注意的不完全分解、预处理、共轭斜量法.

## §1 斜 量 法

设  $A = (a_{ij})$  是  $n \times n$  实对称、正定矩阵. 考虑线性方程组

$$Ax = b \quad (1)$$

的求解问题. 这里  $x = (x_1, x_2, \dots, x_n)^T$  是未知向量,  $b = (b_1, b_2, \dots, b_n)^T$  是已知的向量.

(1) 的求解问题, 等价于下列泛函的求极小问题:

$$F(x) = (Ax, x) - 2(b, x)$$

或

$$F(x) = x^T A x - 2b^T x, \quad (2)$$

即使(2)达到极小的向量  $\tilde{x}$  为(1)的解, 反之, (1)的解是使(2)达到极小的向量.

实际上, 因为  $A$  正定, 故  $A^{-1}$  存在, 记  $\tilde{x} = A^{-1}b$ , 它是(1)的解, 由

$$\begin{aligned}
 F(x) &= (Ax, x) - 2(b, x) = (Ax, x) - 2(A\tilde{x}, x) \\
 &= (A(x - \tilde{x}), (x - \tilde{x})) - (A\tilde{x}, \tilde{x}) \\
 &\geq - (A\tilde{x}, \tilde{x}) = F(\tilde{x})
 \end{aligned}$$

且因为  $A$  是正定的, 故只有当  $x - \tilde{x} = 0$  时, 才能使上式中的  $\geq$  成为等号, 否则是  $>$  号. 这就证明了(1)的求解问题, 等价于(2)的求极小问题.

(2)的求极小问题, 等价于

$$F_1(x) = (A(x - \tilde{x}), (x - \tilde{x})) \quad (3)$$

或

$$F_1(x) = (A^{-1}(Ax - b), (Ax - b)) \quad (4)$$

的求极小问题. 这是因为  $F_1(x) = F(x) + (A\tilde{x}, \tilde{x})$ , 而  $(A\tilde{x}, \tilde{x})$  是常数.

(3)或(4)与(2)的主要差别在于(2)表示式中没有未知量, 因此给定一个  $x$ , 可以算出  $F(x)$ , 而(3), (4)都有未知量出现, 当不知道  $\tilde{x}$  或  $A^{-1}$  时, 给定一个  $x$ , 算不出  $F_1(x)$ . 但当不需要计算泛函值时, 用  $F_1(x)$  更加直观.

一种求(1)解的想法是先取一个初始向量  $x_0$ , 按某种规则求出一个向量  $x_1$ , 使得  $F_1(x_1) < F_1(x_0)$ , 然后再从  $x_1$ , 按上述规则求出一个向量  $x_2$ , 使得  $F_1(x_2) < F_1(x_1)$ , 依此类推, 得到一个向量序列  $x_0, x_1, x_2, \dots$ , 并且希望这个序列能收敛到解  $\tilde{x}$ . 这是一种很一般的想法, 很多方法都是基于这样的想法. 当然所得到的向量序列能不能收敛到  $\tilde{x}$ , 收敛速度如何, 都依赖于按什么样的规则从  $x_i$  确定  $x_{i+1}$ .

斜量法也是一种具体实现这种想法的方法. 它的规则是: 从  $x_0$ , 有  $Ax_0 - b = r_0$ , 如果  $r_0 = 0$ , 那么  $x_0$  即为(1)的解; 如果  $r_0 \neq 0$ , 那么令  $x = x_0 + \alpha r_0$ , 当  $\alpha$  变动时, 表示一条过  $x_0$  的直线, 它的方向跟  $r_0$  相同. 在这条直线上找一点  $x_1 =$

$x_0 + \alpha_0 r_0$ , 使得对所有的实数  $\alpha$ ,

$$F_1(x_1) \leq F_1(x_0 + \alpha r_0), \quad (5)$$

也即在这条直线上,  $x_1$  使  $F_1(x)$  达到极小.

现在来看  $\alpha_0$  是什么? 因为

$$\begin{aligned} F_1(x_0 + \alpha r_0) &= (A(x_0 + \alpha r_0 - \tilde{x}), (x_0 + \alpha r_0 - \tilde{x})) \\ &= (A(x_0 - \tilde{x}), (x_0 - \tilde{x})) \\ &\quad + 2\alpha(A(x_0 - \tilde{x}), r_0) + \alpha^2(Ar_0, r_0), \end{aligned}$$

$$\text{有 } \frac{\partial F_1(x_0 + \alpha r_0)}{\partial \alpha} = 2(A(x_0 - \tilde{x}), r_0) + 2\alpha(Ar_0, r_0),$$

$$\frac{\partial^2 F_1}{\partial \alpha^2} = 2(Ar_0, r_0) > 0,$$

因此取  $\frac{\partial F_1}{\partial \alpha} = 0$  得到的  $\alpha$  即为  $\alpha_0$ , 即

$$\begin{aligned} \alpha_0 &= -(A(x_0 - \tilde{x}), r_0) / (Ar_0, r_0) \\ &= -(r_0, r_0) / (Ar_0, r_0), \end{aligned}$$

从而  $x_1 = x_0 - [(r_0, r_0) / (Ar_0, r_0)] r_0$ , 而

$$\begin{aligned} F_1(x_1) &= (A(x_0 - \tilde{x}), (x_0 - \tilde{x})) \\ &\quad - 2[(r_0, r_0) / (Ar_0, r_0)](r_0, r_0) \\ &\quad + [(r_0, r_0) / (Ar_0, r_0)]^2 (Ar_0, r_0) \\ &= F_1(x_0) - (r_0, r_0)^2 / (Ar_0, r_0) < F_1(x_0). \end{aligned}$$

上面的  $\alpha_0$  是通过  $\frac{\partial F_1}{\partial \alpha} = 0$  而导出的, 也可通过另一种思

想来导出. 对于  $R^n$  引进一种新的内积:  $[x, y] = (Ax, y)$ , 因为  $A$  是对称正定的, 因此它满足内积的条件 (参见 [9]) 从而有一种新的范数  $|x| = \sqrt{(Ax, x)} = \sqrt{[x, x]}$ . 于是

$$F_1(x) = |x - \tilde{x}|^2,$$

$$F_1(x_0 + \alpha r_0) = |x_0 - \tilde{x} + \alpha r_0|^2,$$

求  $x_1$  的问题, 相当于向量  $x_0 - \tilde{x}$  减去  $r_0$  的某一倍数后, 使

余量最小. 自然  $\alpha_0$  使得  $\mathbf{x}_0 - \tilde{\mathbf{x}} + \alpha_0 \mathbf{r}_0$  正交于  $\mathbf{r}_0$ . 即  $\alpha_0$  满足

$$[\mathbf{x}_0 - \tilde{\mathbf{x}} + \alpha_0 \mathbf{r}_0, \mathbf{r}_0] = 0, \\ \alpha_0 = -\frac{[\mathbf{x}_0 - \tilde{\mathbf{x}}, \mathbf{r}_0]}{[\mathbf{r}_0, \mathbf{r}_0]} = -\frac{(\mathbf{r}_0, \mathbf{r}_0)}{(A\mathbf{r}_0, \mathbf{r}_0)} \quad (6)$$

时, 可使余量  $|\mathbf{x}_0 + \alpha_0 \mathbf{r}_0 - \tilde{\mathbf{x}}|$  最小.

这种由正交性导出极值性的办法是有普遍意义的.

**定理 1.1**  $R^n$  上定义的任意一种内积  $((\cdot, \cdot))$ ,  $\mathbf{y}$  是  $R^n$  中的任意一个向量,  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_k$  是  $R^n$  中  $k$  个线性无关的向量, 则  $k$  个常数  $\alpha_1^0, \alpha_2^0, \dots, \alpha_k^0$  使

$$((\mathbf{y} + \alpha_1^0 \mathbf{f}_1 + \dots + \alpha_k^0 \mathbf{f}_k, \mathbf{y} + \alpha_1^0 \mathbf{f}_1 + \dots + \alpha_k^0 \mathbf{f}_k)) \\ < ((\mathbf{y} + \alpha_1 \mathbf{f}_1 + \dots + \alpha_k \mathbf{f}_k, \mathbf{y} + \alpha_1 \mathbf{f}_1 + \dots + \alpha_k \mathbf{f}_k)), \quad (7)$$

对任意实数  $\alpha_1, \alpha_2, \dots, \alpha_k$  只要  $(\alpha_1, \alpha_2, \dots, \alpha_k) \neq (\alpha_1^0, \alpha_2^0, \dots, \alpha_k^0)$ , 成立的充分必要条件是:  $\alpha_1^0, \alpha_2^0, \dots, \alpha_k^0$  满足

$$((\mathbf{y} + \alpha_i^0 \mathbf{f}_1 + \dots + \alpha_k^0 \mathbf{f}_k, \mathbf{f}_i)) = 0, \quad i = 1, 2, \dots, k. \quad (8)$$

**证明** 若(8)成立, 记  $\alpha_i = \alpha_i^0 + \Delta\alpha_i$ , 则

$$\begin{aligned} & \left( \left( \mathbf{y} + \sum_{i=1}^k \alpha_i \mathbf{f}_i, \mathbf{y} + \sum_{i=1}^k \alpha_i \mathbf{f}_i \right) \right) \\ &= \left( \left( \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i + \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i, \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i + \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i \right) \right) \\ &= \left( \left( \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i, \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i \right) \right) \\ &\quad + 2 \left( \left( \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i, \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i \right) \right) \\ &\quad + \left( \left( \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i, \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i \right) \right) \end{aligned}$$

中间项  $\left( \left( \mathbf{y} + \sum_{i=1}^k \alpha_i^0 \mathbf{f}_i, \sum_{i=1}^k \Delta\alpha_i \mathbf{f}_i \right) \right)$  为 0, 而第 3 项因为  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_k$  线性无关, 除非所有  $\Delta\alpha_i (i=1, 2, \dots, k)$  为 0, 否则总是正的, 因此(7)成立.

反之若(7)成立, 取

$$(\alpha_1, \alpha_2, \dots, \alpha_k) = (\alpha_1^0, \dots, \alpha_{i-1}^0, \alpha_i^0 + \Delta\alpha_i, \alpha_{i+1}^0, \dots, \alpha_k^0),$$

于是

$$\begin{aligned} & \left( \left( \mathbf{y} + \sum_{l=1}^k \alpha_l \mathbf{f}_l, \mathbf{y} + \sum_{l=1}^k \alpha_l \mathbf{f}_l \right) \right) \\ &= \left( \left( \mathbf{y} + \sum_{l=1}^k \alpha_l^0 \mathbf{f}_l, \mathbf{y} + \sum_{l=1}^k \alpha_l^0 \mathbf{f}_l \right) \right) \\ & \quad + 2\Delta\alpha_i \left( \left( \mathbf{y} + \sum_{l=1}^k \alpha_l^0 \mathbf{f}_l, \mathbf{f}_i \right) \right) + (\Delta\alpha_i)^2 ((\mathbf{f}_i, \mathbf{f}_i)), \end{aligned}$$

由(7)知  $\Delta\alpha_i = 0$  使上式达到极小, 故必须  $2\Delta\alpha_i$  的系数

$$\left( \left( \mathbf{y} + \sum_{l=1}^k \alpha_l^0 \mathbf{f}_l, \mathbf{f}_i \right) \right) = 0,$$

否则因为  $((\mathbf{f}_i, \mathbf{f}_i)) > 0$ , 取

$$\Delta\alpha_i = - \left( \left( \mathbf{y} + \sum_{l=1}^k \alpha_l^0 \mathbf{f}_l, \mathbf{f}_i \right) \right) / ((\mathbf{f}_i, \mathbf{f}_i)) \neq 0,$$

将使该式达到极小, 跟假设(7)成立矛盾.

因为所取的  $i$  是任意的, 因此(8)成立. 定理证毕.

再回到方程(1)的求解问题上来. 求得了  $\mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{r}_0$  后, 可以构造出  $\mathbf{r}_1 = A\mathbf{x}_1 - \mathbf{b} = \mathbf{r}_0 + \alpha_0 A\mathbf{r}_0$ , 于是可以在直线  $\mathbf{x} = \mathbf{x}_1 + \alpha \mathbf{r}_1$  上求一点  $\mathbf{x}_2 = \mathbf{x}_1 + \alpha_1 \mathbf{r}_1$  使得

$$F_1(\mathbf{x}_1 + \alpha_1 \mathbf{r}_1) < F_1(\mathbf{x}_1 + \alpha \mathbf{r}_1),$$

对任意  $\alpha \neq \alpha_1$  的实数成立. 这样的

$$\alpha_1 = -(\mathbf{r}_1, \mathbf{r}_1) / (A\mathbf{r}_1, \mathbf{r}_1),$$

依此类推, 有计算程式:

给定  $\mathbf{x}_0$ ,

$$\begin{cases} \mathbf{r}_k = A\mathbf{x}_k - \mathbf{b}, \\ \alpha_k = -(\mathbf{r}_k, \mathbf{r}_k) / (A\mathbf{r}_k, \mathbf{r}_k), \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k, \quad k=0, 1, 2, \dots, \end{cases} \quad (9)$$

从(9)构造出序列  $\{x_k\}$  的方法就称为斜量法, 因为极小是在  $r_k$  的方向上取的, 而

$$r_k = \frac{1}{2} \operatorname{grad} F_1(x) \Big|_{x=x_k},$$

因此称为斜量法.

斜量又称梯度 (gradient), 它的几何意义是使  $F_1(x)$  在某点邻近变化最快的方向, 因此从  $F_1(x_0 + \alpha r_0)$  求极小, 希望比其他方向  $p$  上求  $F_1(x_0 + \alpha p)$  极小, 使得  $F_1(x)$  下降得更快一点.

下面定理回答了从(9)导出的序列  $\{x_k\}$  的收敛性问题.

**定理 1.2** 设  $A$  的特征值为  $\lambda_1, \lambda_2, \dots, \lambda_n$ ,

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n,$$

则

$$\|x_k - \tilde{x}\| \leq \sqrt{\frac{\lambda_n}{\lambda_1}} \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^k \|x_0 - \tilde{x}\|. \quad (10)$$

为了证明(10)式成立, 先要证明两个引理. 设矩阵  $A$  的特征值  $\lambda_i$  对应的单位特征向量为  $y_i$ , 并且  $y_1, y_2, \dots, y_n$  是一组标准正交向量组.

**引理 1.1** 在定理 1.2 的假设下, 成立下列不等式

$$\lambda_1(x, x) \leq (Ax, x) \leq \lambda_n(x, x).$$

**证明** 令  $x = \sum_{j=1}^n \beta_j y_j$ , 于是

$$(Ax, x) = \sum_{j=1}^n \lambda_j \beta_j^2,$$

$$\lambda_1(x, x) = \sum_{j=1}^n \lambda_1 \beta_j^2,$$

$$\lambda_n(x, x) = \sum_{j=1}^n \lambda_n \beta_j^2,$$

故  $\lambda_1(x, x) \leq (Ax, x) \leq \lambda_n(x, x)$ . 证毕.

引理 1.2 设  $\varphi(\lambda)$  是一个  $\lambda$  的多项式, 则

$$[\varphi(A)x, \varphi(A)x] \leq \max_i \varphi(\lambda_i)^2 [x, x].$$

证明 令  $x = \sum_{j=1}^n \beta_j y_j$ ,

$$\varphi(A)x = \sum_{j=1}^n \varphi(\lambda_j) \beta_j y_j,$$

$$[\varphi(A)x, \varphi(A)x] = \sum_{j=1}^n \lambda_j \varphi(\lambda_j)^2 \beta_j^2$$

$$\leq \max_i \varphi(\lambda_i)^2 \sum_{j=1}^n \lambda_j \beta_j^2 = \max_i \varphi(\lambda_i)^2 [x, x]. \text{ 证毕.}$$

现在来证明定理 1.2. 若  $r_0 = 0$ , (10) 式自然成立, 假设  $r_0 \neq 0$ , 由  $F_1(x_1) \leq F_1(x_0 + \alpha r_0)$  知

$$\begin{aligned} [x_1 - \tilde{x}, x_1 - \tilde{x}] &\leq [x_0 + \alpha r_0 - \tilde{x}, x_0 + \alpha r_0 - \tilde{x}] \\ &= [(I + \alpha A)(x_0 - \tilde{x}), (I + \alpha A)(x_0 - \tilde{x})], \end{aligned}$$

对于多项式  $1 + \alpha\lambda$ , 应用引理 1.2, 得到

$$\begin{aligned} [x_1 - \tilde{x}, x_1 - \tilde{x}] &\leq \max_i (1 + \alpha\lambda_i)^2 [x_0 - \tilde{x}, x_0 - \tilde{x}] \\ &\leq \max_{\lambda_1 \leq \lambda \leq \lambda_n} (1 + \alpha\lambda)^2 [x_0 - \tilde{x}, x_0 - \tilde{x}], \end{aligned}$$

这一不等式对任何实数  $\alpha$  都成立的, 因此对特别选取的  $\alpha =$

$$-\frac{2}{\lambda_1 + \lambda_n}, \text{ 也成立不等式}$$

$$[x_1 - \tilde{x}, x_1 - \tilde{x}]$$

$$\leq \max_{\lambda_1 \leq \lambda \leq \lambda_n} \left(1 - \frac{2\lambda}{\lambda_1 + \lambda_n}\right)^2 [x_0 - \tilde{x}, x_0 - \tilde{x}],$$

但

$$\max_{\lambda_1 \leq \lambda \leq \lambda_n} \left(1 - \frac{2\lambda}{\lambda_1 + \lambda_n}\right)^2 = \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^2,$$

所以

$$[x_1 - \tilde{x}, x_1 - \tilde{x}] \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^2 [x_0 - \tilde{x}, x_0 - \tilde{x}],$$



这样的关系, 对于  $x_k, x_{k-1}$  之间也可以同样证明成立的, 即

$$[x_k - \tilde{x}, x_k - \tilde{x}] \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 [x_{k-1} - \tilde{x}, x_{k-1} - \tilde{x}],$$

而  $k$  又是任意自然数, 故

$$[x_k - \tilde{x}, x_k - \tilde{x}] \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^{2k} [x_0 - \tilde{x}, x_0 - \tilde{x}],$$

再利用引理 1.1 就可获得估计式 (10). 证毕.

在证明中取  $\alpha = -\frac{2}{\lambda_1 + \lambda_n}$ , 它是下列极大极小问题的解,

$$\min_{\alpha} \max_{\lambda_1 < \lambda < \lambda_n} |1 + \alpha \lambda|$$

即对任何实数  $\alpha$  有

$$\max_{\lambda_1 < \lambda < \lambda_n} \left| 1 - \frac{2}{\lambda_1 + \lambda_n} \lambda \right| \leq \max_{\lambda_1 < \lambda < \lambda_n} |1 + \alpha \lambda|, \quad (11)$$

等式只有在  $\alpha = -2/(\lambda_1 + \lambda_n)$  时达到, 这可从图 1 知道. 由 (11) 式可知取  $\alpha = -2/(\lambda_1 + \lambda_n)$  可以得到比较小的误差估计, 因此在证明中取这一个值.

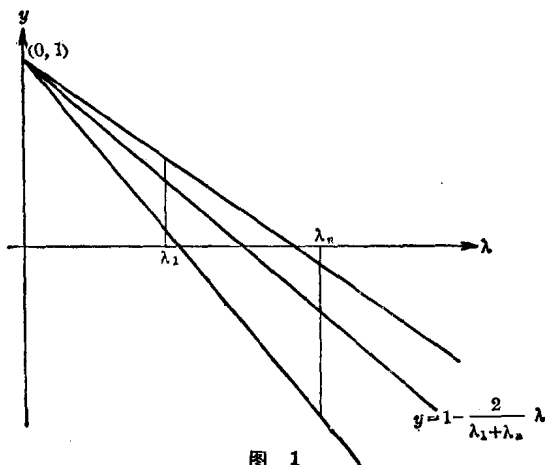


图 1

因为  $0 < \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} < 1$ , 从定理 1.2 知斜量法所得的序列  $\{x_k\}$  是收敛的。

称  $p = \|A\|A^{-1}$  为系数矩阵  $A$  的条件数。当  $A$  是对称正定时  $p = \lambda_n/\lambda_1$ , 因此  $\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} = \frac{p-1}{p+1}$ , 从而可知当  $p$  很大时  $\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}$  接近于 1。说明矩阵  $A$  的条件数很大时 (此时称矩阵  $A$  是病态的), (10) 式右端收敛于 0 是很慢的, 实际计算的经验也是这样, 此时  $x_k$  收敛于  $\tilde{x}$  的速度也是很慢。

为了对系数矩阵  $A$  的条件数有个数量的概念, 我们用差分方法求解正方形  $0 \leq x, y \leq 1$  上的 Poisson 方程第一类边值问题为例, 若差分步长为  $h=1/m$ , 则系数矩阵为  $(m-1)^2$  阶, 它的  $(m-1)^2$  个特征值为

$$\lambda_{p,q} = 1 - \frac{1}{2} (\cos p\pi h + \cos q\pi h) \quad p, q = 1, 2, \dots, m-1,$$

故 
$$\lambda_1 = 1 - \cos h\pi = 2 \sin^2 \frac{h\pi}{2},$$

$$\lambda_n = 1 - \cos(m-1)h\pi = 2 \sin^2 \frac{(m-1)\pi}{2m},$$

于是 
$$p = \frac{\lambda_n}{\lambda_1} \approx 4 / (h\pi)^2,$$

$$h = 1/10, \quad p \approx 39.87, \quad \frac{p-1}{p+1} = 0.95,$$

$$h = 1/50, \quad p \approx 1000, \quad \frac{p-1}{p+1} = 0.998,$$

$$h = 1/100, \quad p \approx 4 \times 10^3, \quad \frac{p-1}{p+1} = 0.9995,$$

当  $h = 1/10$ ,  $\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} = 0.95$ , 如果要

$$\left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^k \leq 10^{-3},$$

$$k \geq -3/\log 0.95 = 134.65,$$

这告诉我们即使对于  $h=1/10$ , 对 81 阶的系数矩阵, 要使

$$\|x_k - \tilde{x}\| \leq 10^{-3} \sqrt{\frac{\lambda_n}{\lambda_1}} \|x_0 - \tilde{x}\|,$$

也要迭代 135 次; 这就表示斜量法的收敛速度显得太慢了, 因此在目前求解线性代数方程组时, 不再采用斜量法. 不过有关斜量法的一些想法对于构造新的方法却仍然是有价值的.

## §2 多步斜量法

上节所介绍的斜量法, 它从  $x_0$  求  $x_1$ , 是在直线  $x_0 + \alpha r_0$  上找的. 同样可以考虑, 从  $x_0$  求  $x_1$ , 在一个平面  $x_0 + \alpha p + \beta q$  上找, 也即  $x_1 = x_0 + \alpha_0 p + \beta_0 q$ , 并且不等式

$$F_1(x_1) \leq F_1(x_0 + \alpha p + \beta q)$$

对一切实数  $\alpha, \beta$  成立. 利用定理 1.1 知道这样的  $\alpha_0, \beta_0$  由下列方程组所确定:

$$\begin{cases} [x_0 + \alpha_0 p + \beta_0 q - \tilde{x}, p] = 0, \\ [x_0 + \alpha_0 p + \beta_0 q - \tilde{x}, q] = 0. \end{cases}$$

只要  $p, q$  线性无关, 上述方程唯一确定  $\alpha_0$  和  $\beta_0$ .

特别取  $p = r_0, q$  为另一与  $r_0$  线性无关的向量时,  $x_1 = x_0 + \alpha_0 r_0 + \beta_0 q$ , 有可能比斜量法中所确定的  $x_1$  要好, 不可能比它差.

另外在 §1 估计斜量法的收敛速度时, 用到等式

$$x_0 + \alpha r_0 - \tilde{x} = (1 + \alpha A)(x_0 - \tilde{x}),$$

从而把问题归结为考虑多项式  $1 + \alpha \lambda$  在区间  $[\lambda_1, \lambda_n]$  上的极大极小问题

$$\min_{\alpha} \max_{\lambda_1 \leq \lambda \leq \lambda_n} |1 + \alpha \lambda|.$$

如果我们仍要把在平面  $\mathbf{x}_0 + \alpha \mathbf{p} + \beta \mathbf{q}$  上找  $\mathbf{x}_1$  的问题, 归结为多项式的极大极小问题, 自然会想到取  $\mathbf{q} = A\mathbf{r}_0$ , 这样

$$\mathbf{x}_0 + \alpha \mathbf{p} + \beta \mathbf{q} = \mathbf{x}_0 + \alpha_0 \mathbf{r}_0 + \beta A\mathbf{r}_0,$$

而

$$\begin{aligned} \mathbf{x}_0 + \alpha \mathbf{r}_0 + \beta A\mathbf{r}_0 - \tilde{\mathbf{x}} &= \mathbf{x}_0 - \tilde{\mathbf{x}} + \alpha \mathbf{r}_0 + \beta A\mathbf{r}_0 \\ &= (I + \alpha A + \beta A^2)(\mathbf{x}_0 - \tilde{\mathbf{x}}). \end{aligned}$$

如果  $\mathbf{r}_0$  与  $A\mathbf{r}_0$  线性无关, 就可以唯一确定

$$\mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{r}_0 + \beta_0 A\mathbf{r}_0,$$

如果  $\mathbf{r}_0$  与  $A\mathbf{r}_0$  线性相关, 我们来分析一下, 会发生什么? 假定  $\mathbf{r}_0 \neq 0$ , 此时必有  $a \neq 0$ , 而

$$A\mathbf{r}_0 = a\mathbf{r}_0,$$

因为  $A^{-1}$  存在, 故

$$\mathbf{r}_0 = a(\mathbf{x}_0 - \tilde{\mathbf{x}}),$$

或

$$\tilde{\mathbf{x}} = \mathbf{x}_0 - \frac{1}{a} \mathbf{r}_0,$$

说明解向量  $\tilde{\mathbf{x}}$  在直线  $\mathbf{x}_0 + \alpha \mathbf{r}_0$  上, 因此  $F_1(\mathbf{x}_0 + \alpha \mathbf{r}_0 + \beta A\mathbf{r}_0)$  上的极小一定在  $\tilde{\mathbf{x}} = \mathbf{x}_0 - \frac{1}{a} \mathbf{r}_0$  上达到.  $\mathbf{x}_1$  即为解  $\tilde{\mathbf{x}}$ . 不管  $\mathbf{r}_0$  与  $A\mathbf{r}_0$  是否线性无关,  $\alpha_0, \beta_0$  都可由

$$\begin{cases} [\mathbf{x}_0 + \alpha_0 \mathbf{r}_0 + \beta_0 A\mathbf{r}_0 - \tilde{\mathbf{x}}, \mathbf{r}_0] = 0, \\ [\mathbf{x}_0 + \alpha_0 \mathbf{r}_0 + \beta_0 A\mathbf{r}_0 - \tilde{\mathbf{x}}, A\mathbf{r}_0] = 0 \end{cases}$$

确定. 一般的计算程式为:

取定一个  $\mathbf{x}_0$ , 计算  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b}$ , 计算  $A\mathbf{r}_0$ , 求解  $\alpha_k, \beta_k$ :

$$\begin{cases} [\mathbf{r}_k, \mathbf{r}_k] \alpha_k + [A\mathbf{r}_k, \mathbf{r}_k] \beta_k = -[\mathbf{x}_k - \tilde{\mathbf{x}}, \mathbf{r}_k], \\ [\mathbf{r}_k, A\mathbf{r}_k] \alpha_k + [A\mathbf{r}_k, A\mathbf{r}_k] \beta_k = -[\mathbf{x}_k - \tilde{\mathbf{x}}, A\mathbf{r}_k], \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k + \beta_k A\mathbf{r}_k, \\ \mathbf{r}_{k+1} = A\mathbf{x}_{k+1} - \mathbf{b}, \quad k=0, 1, 2, \dots, \end{cases} \quad (12)$$



$\alpha_k^{(2)}, \dots, \alpha_k^{(l)}$ :

$$\left\{ \begin{aligned} & [\mathbf{r}_k, \mathbf{r}_k] \alpha_k^{(1)} + [A\mathbf{r}_k, \mathbf{r}_k] \alpha_k^{(2)} + \dots + [A^{l-1}\mathbf{r}_k, \mathbf{r}_k] \alpha_k^{(l)} \\ & \quad = -[\mathbf{x}_k - \tilde{\mathbf{x}}, \mathbf{r}_k], \\ & [\mathbf{r}_k, A\mathbf{r}_k] \alpha_k^{(1)} + [A\mathbf{r}_k, A\mathbf{r}_k] \alpha_k^{(2)} + \dots + [A^{l-1}\mathbf{r}_k, A\mathbf{r}_k] \alpha_k^{(l)} \\ & \quad = -[\mathbf{x}_k - \tilde{\mathbf{x}}, A\mathbf{r}_k], \\ & \dots\dots\dots \\ & [\mathbf{r}_k, A^{l-1}\mathbf{r}_k] \alpha_k^{(1)} + [A\mathbf{r}_k, A^{l-1}\mathbf{r}_k] \alpha_k^{(2)} + \dots \\ & \quad + [A^{l-1}\mathbf{r}_k, A^{l-1}\mathbf{r}_k] \alpha_k^{(l)} = -[\mathbf{x}_k - \tilde{\mathbf{x}}, A^{l-1}\mathbf{r}_k], \\ & \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k^{(1)}\mathbf{r}_k + \alpha_k^{(2)}A\mathbf{r}_k + \dots + \alpha_k^{(l)}A^{l-1}\mathbf{r}_k, \\ & \mathbf{r}_{k+1} = A\mathbf{x}_{k+1} - \mathbf{b}, \quad k=0, 1, 2, \dots \end{aligned} \right. \quad (13)$$

记

$$E_l = \min_{\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(l)}} \max_{\lambda_1 \leq \lambda < \lambda_n} |1 + \alpha^{(1)}\lambda + \alpha^{(2)}\lambda^2 + \dots + \alpha^{(l)}\lambda^l|,$$

同样对  $l$  步斜量法有收敛速度的估计

$$\begin{aligned} |\mathbf{x}_1 - \tilde{\mathbf{x}}| &\leq E_l |\mathbf{x}_0 - \tilde{\mathbf{x}}|, \\ \|\mathbf{x}_k - \tilde{\mathbf{x}}\| &\leq \sqrt{\frac{\lambda_n}{\lambda_1}} E_l^k \|\mathbf{x}_0 - \tilde{\mathbf{x}}\|. \end{aligned} \quad (14)$$

在 §1 已经知道

$$E_1 = \min_{\alpha} \max_{\lambda_1 \leq \lambda < \lambda_n} |1 + \alpha\lambda| = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1},$$

对于  $l > 1$ ,  $E_l$  究竟是怎样一个量? 这是一个需要进一步回答的问题.  $E_l$  是从型如  $1 + \alpha^{(1)}\lambda + \alpha^{(2)}\lambda^2 + \dots + \alpha^{(l)}\lambda^l$  的多项式全体中考虑极大极小得到的量, 也即是从所有  $l$  次多项式  $q_l(\lambda)$  中, 满足  $q_l(0) = 1$  的那些  $l$  次多项式中考虑极大极小得到的量.

下面先介绍一下契贝谢夫 (П. Л. Чебышев) 多项式.

$$T_l(x) = \frac{(x + \sqrt{x^2 - 1})^l + (x - \sqrt{x^2 - 1})^l}{2}$$

是一个首项系数为  $2^{l-1}$  的  $l$  次多项式, 称为  $l$  次契贝谢夫多

项式. 如果把自变量  $x$  的范围限制在  $[-1, 1]$  上, 那么可以作变换  $x = \cos \theta$ ,  $\theta \in [0, \pi]$ ,  $\theta$  与  $x$  之间建立起一一对应.

$$\begin{aligned} T_l(\cos \theta) &= \frac{(\cos \theta + \sqrt{\cos^2 \theta - 1})^l + (\cos \theta - \sqrt{\cos^2 \theta - 1})^l}{2} \\ &= \frac{(\cos \theta + i \sin \theta)^l + (\cos \theta - i \sin \theta)^l}{2} = \cos l\theta, \end{aligned}$$

因此当  $x \in [-1, 1]$  时,

$$T_l(x) = \cos l \arccos x,$$

从而也有  $|T_l(x)| \leq 1$ .

取  $x_k = \cos \frac{k\pi}{l} \in [-1, 1]$ , 有

$$T_l(x_k) = \cos l \arccos \cos \frac{k\pi}{l} = \cos k\pi = (-1)^k,$$

对于  $k=0, 1, 2, \dots, l$ ,  $x_0, x_1, \dots, x_l$  是  $l+1$  个从右到左排列在区间  $[-1, 1]$  上不同的点, 它们轮流使  $T_l(x)$  取正 1 或负 1, 而 1 是  $T_l(x)$  在  $[-1, 1]$  上的最大值, 这样的一组点称为正负交替偏差点. 有  $l+1$  个正负交替偏差点是  $l$  次契贝谢夫多项式的一个重要性质.

今考虑如下  $l$  次多项式

$$P_l(\lambda) = T_l\left(\frac{2\lambda}{\lambda_n - \lambda_1} - \frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1}\right) / T_l\left(-\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1}\right),$$

显然有  $P_l(0) = 1$ , 并且当  $\lambda \in [\lambda_1, \lambda_n]$  时,

$$x = \frac{2\lambda}{\lambda_n - \lambda_1} - \frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} \in [-1, 1],$$

因此当  $\lambda \in [\lambda_1, \lambda_n]$  时

$$|P_l(\lambda)| \leq 1 / \left| T_l\left(-\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1}\right) \right|,$$

将  $[-1, 1]$  上的  $x_k$  与  $[\lambda_1, \lambda_n]$  中的  $\eta_k$  建立起对应, 即

$$x_k = \frac{2\eta_k}{\lambda_n - \lambda_1} - \frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1},$$

或 
$$\eta_k = \left( \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} + \cos \frac{k\pi}{l} \right) \frac{\lambda_n - \lambda_1}{2},$$

有  $P_l(\eta_k) = (-1)^k / T_l \left( -\frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} \right)$ , 这说明  $\eta_k (k=0, 1, \dots, l)$  是  $P_l(t)$  在  $[\lambda_1, \lambda_n]$  中的  $l+1$  个正负交替的偏差点.

**定理 1.3** 
$$E_l = \max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)| = 1 / T_l \left( \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right).$$

**证明** 由  $E_l$  的定义可知:

$$E_l \leq \max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)|,$$

如果  $E_l < \max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)|$ , 那么必有一个  $l$  次多项式  $Q(\lambda)$ ,

且  $Q(0) = 1$ , 使得

$$\max_{\lambda_1 < \lambda < \lambda_n} |Q(\lambda)| < \max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)|, \quad (15)$$

于是  $P_l(\lambda) - Q(\lambda)$  在  $\eta_k (k=0, 1, \dots, l)$  的符号与  $(-1)^k / T_l \left( -\frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} \right)$  相同, 从而  $P_l(\lambda) - Q(\lambda)$  在  $[\lambda_1, \lambda_n]$  中至少有  $l$  个零点, 另外 0 不在区间  $[\lambda_1, \lambda_n]$  中, 而

$$P_l(0) - Q(0) = 1 - 1 = 0,$$

因此 0 也是  $P_l(\lambda) - Q(\lambda)$  的一个零点, 这样  $P_l(\lambda) - Q(\lambda)$  至少有  $l+1$  个零点, 因为这是一个  $l$  次多项式, 故必须

$$P_l(\lambda) - Q(\lambda) \equiv 0,$$

此与 (15) 矛盾, 故知

$$E_l = \max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)|,$$

但 
$$\max_{\lambda_1 < \lambda < \lambda_n} |P_l(\lambda)| = 1 / \left| T_l \left( -\frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} \right) \right|,$$

而对于契贝谢夫多项式有



$$T_l\left(-\frac{\lambda_1+\lambda_n}{\lambda_n-\lambda_1}\right)=(-1)^l T_l\left(\frac{\lambda_1+\lambda_n}{\lambda_n-\lambda_1}\right),$$

故  $E_l=1/T_l\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}\right)$ . 证毕.

按  $T_l(x)$  的定义

$$\begin{aligned} E_l &= 1/T_l\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}\right) \\ &= 2/\left[\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}+\sqrt{\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}\right)^2-1}\right)^l\right. \\ &\quad \left.+\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}-\sqrt{\left(\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}\right)^2-1}\right)^l\right], \end{aligned}$$

$$\text{当 } l=1 \text{ 时 } E_1 = \frac{2}{\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1} + \frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}} = \frac{\lambda_n-\lambda_1}{\lambda_n+\lambda_1},$$

与 §1 的结果相符.

定理 1.3 没有说明在  $1+\alpha^{(1)}\lambda+\cdots+\alpha^{(l)}\lambda^l$  型的多项式中, 是否只有一个  $P_l(\lambda)$ , 使得

$$\max_{\lambda_1 < \lambda < \lambda_n} |1+\alpha^{(1)}\lambda+\alpha^{(2)}\lambda^2+\cdots+\alpha^{(l)}\lambda^l| = E_l.$$

这个问题的回答是肯定的.

**定理 1.4** 如果  $l$  次多项式  $Q(\lambda)$ ,  $Q(0)=1$ , 满足如下条件

$$\max_{\lambda_1 < \lambda < \lambda_n} |Q(\lambda)| = E_l,$$

则  $Q(\lambda) \equiv P_l(\lambda)$ .

**证明** 由定理 1.3 知  $P_l(\lambda)-Q(\lambda)$  在  $\eta_k$  上或者同号于  $(-1)^k/T_l\left(-\frac{\lambda_n+\lambda_1}{\lambda_n-\lambda_1}\right)$ , 或者为 0. 对于  $\eta_k (k=1, 2, \dots, l-1)$  有  $\eta_k \in (\lambda_1, \lambda_n)$ , 因此当在这些内部的  $\eta_k$  如果

$$P_l(\eta_k)-Q(\eta_k)=0,$$

必有  $P_l'(\eta_k) - Q'(\eta_k) = 0$ , 这是因为  $\eta_k$  是  $P_l(\lambda)$  的极大(小)值, 也是  $Q(\lambda)$  的极大(小)值. 这时,  $\eta_k$  是  $P_l(\lambda) - Q(\lambda)$  的重根. 由此可知  $P_l(\lambda) - Q(\lambda)$  在  $[\lambda_1, \lambda_n]$  中至少有  $l$  个零点. 另外再由  $P_l(0) - Q(0) = 0$ , 知  $0$  也是一个零点, 从而可知  $P_l(\lambda) - Q(\lambda)$  至少有  $l+1$  个零点, 故

$$Q(\lambda) \equiv P_l(\lambda). \text{ 证毕.}$$

现在来比较斜量法与多步斜量法的收敛速度. 在计算中, 主要的计算量是化在矩阵与向量的乘法上, 我们以矩阵与向量乘法一次作为一个单位.

对斜量法来说从  $\mathbf{x}_0$  求  $\mathbf{x}_1$  要计算二次矩阵和向量的乘法, 即  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b}$ ,  $A\mathbf{r}_0$ , 但对以后从  $\mathbf{x}_k$  求  $\mathbf{x}_{k+1}$ , 每步只要计算一次矩阵和向量的乘法  $A\mathbf{r}_k$ .

对于二步斜量法来说, 从  $\mathbf{x}_0$  求  $\mathbf{x}_1$  要计算三次矩阵和向量的乘法,  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b}$ ,  $A\mathbf{r}_0$ ,  $A^2\mathbf{r}_0$ , 但对以后从  $\mathbf{x}_k$  求  $\mathbf{x}_{k+1}$  每步只要二次矩阵和向量的乘法  $A\mathbf{r}_k$ ,  $A^2\mathbf{r}_k$ .

如果我们称斜量法为一步斜量法, 那么不算从  $\mathbf{x}_0$  求  $\mathbf{x}_1$ , 一般地从  $\mathbf{x}_k$  求  $\mathbf{x}_{k+1}$ ,  $l$  步斜量法需要计  $l$  次矩阵和向量的乘法. 因此  $l$  步斜量迭代一次, 计算量相当于斜量法迭代  $l$  次. 因此我们要比较  $E_l$  与  $E_1^l$  究竟那个小, 才能知道那个收敛速度快. 同样各种多步斜量法之间也要作比较,  $l$  步斜量法与  $k$  步斜量法究竟那个快, 就要比较  $(E_l)^{1/l}$  与  $(E_k)^{1/k}$  究竟那个小. 下面定理回答了这个问题.

$$\text{定理 1.5 设 } T_s(\sigma) = \frac{(\sigma + \sqrt{\sigma^2 - 1})^s + (\sigma - \sqrt{\sigma^2 - 1})^s}{2}$$

是  $s$  次契贝谢夫多项式. 若  $\sigma > 1$ ,  $k > l$ , 则

$$(T_k(\sigma))^{1/k} > (T_l(\sigma))^{1/l}.$$

证明 令  $\alpha = \sigma + \sqrt{\sigma^2 - 1}$ ,  $\beta = \sigma - \sqrt{\sigma^2 - 1}$ , 于是

$$\theta = \beta/\alpha < 1,$$

而  $T_s(\sigma) = \alpha^s \frac{1+\theta^s}{2}$ ,  $(T_s(\sigma))^{1/s} = \alpha \left( \frac{1+\theta^s}{2} \right)^{1/s}$ , 我们下面证明  $(T_s(\sigma))^{1/s}$  是  $s$  的单调上升函数, 为此只要证明

$$f(s) = \left( \frac{1+\theta^s}{2} \right)^{1/s}$$

是  $s$  的单调上升函数.

实际上

$$\begin{aligned} f'(s) &= f(s) \frac{d}{ds} \left[ \frac{1}{s} \ln \frac{1+\theta^s}{2} \right] \\ &= f(s) \left\{ -\frac{1}{s^2} \ln \frac{1+\theta^s}{2} + \frac{1}{s} \frac{2}{1+\theta^s} \times \frac{1}{2} \theta^s \ln \theta \right\} \\ &= f(s) \left\{ -\frac{1}{s^2} \ln \frac{1+\theta^s}{2} + \frac{s\theta^s \ln \theta}{s^2(1+\theta^s)} \right\} \\ &= f(s) \frac{(1+\theta^s) \ln 2 - (1+\theta^s) \ln(1+\theta^s) + \theta^s \ln \theta^s}{s^2(1+\theta^s)}. \end{aligned}$$

为了判明  $f'(s)$  的符号, 我们只要研究上述等式右边的分子符号, 如果记  $\eta = \theta^s$ , 于是知  $\eta \in [0, 1)$ , 上述分子为

$$g(\eta) = (1+\eta) \ln 2 - (1+\eta) \ln(1+\eta) + \eta \ln \eta,$$

$$g'(\eta) = \ln 2 - \ln(1+\eta) - 1 + \ln \eta + 1$$

$$= \ln \frac{2\eta}{1+\eta} < 0, \quad \text{当 } \eta \in (0, 1),$$

$$g'(1) = 0,$$

再由  $g(0) = \ln 2$ ,  $g(1) = 0$ , 利用微分学中值定理就知  $g(\eta) > 0$  当  $\eta \in [0, 1)$ . 再利用  $f(s) > 0$ ,  $s^2(1+\theta^s) > 0$  得到

$$f'(s) > 0. \text{ 证毕.}$$

**推论** 当  $k > l$ ,  $(E_k)^{\frac{1}{k}} < (E_l)^{\frac{1}{l}}$ .

**证明**  $\sigma = \frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} > 1,$

而  $E_s = 1/T_s(\sigma)$ , 利用定理 1.5 即可证明推论.

从定理 1.5 知道, 多步斜量法比斜量法要快, 而且步数越高越快, 但是步数越高 (13) 中的方程组阶数也越高, 需要更多的存贮单元, 需要更多的计算量求解也越麻烦, 下节介绍的共轭斜量法, 避免了解 (13) 中的那样方程组。

### § 3 共轭斜量法

给定向量  $\mathbf{x}_0$ ,  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b}$ , 可以得到序列

$$\mathbf{r}_0, A\mathbf{r}_0, \dots, A^l\mathbf{r}_0, \quad l < n,$$

假定这组向量是线性无关的, 记子空间  $m_l = \{\mathbf{r}_0, A\mathbf{r}_0, \dots, A^l\mathbf{r}_0\}$ , 将这组向量按内积  $[\mathbf{x}, \mathbf{y}]$  正交化, 这样的正交化也称为  $A$ -正交化, 于是得到一组  $A$ -正交化的序列

$$\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_l,$$

假定正交化过程是如下那样次序进行的:

$$\begin{cases} \mathbf{q}_0 = \mathbf{r}_0, \mathbf{q}_1 = A\mathbf{r}_0 + \sigma_0^{(0)}\mathbf{q}_0, \\ \mathbf{q}_2 = A^2\mathbf{r}_0 + \sigma_0^{(1)}\mathbf{q}_0 + \sigma_1^{(1)}\mathbf{q}_1, \\ \text{一般} \\ \mathbf{q}_i = A^i\mathbf{r}_0 + \sigma_0^{(i-1)}\mathbf{q}_0 + \sigma_1^{(i-1)}\mathbf{q}_1 + \dots + \sigma_{i-1}^{(i-1)}\mathbf{q}_{i-1}, \end{cases} \quad (16)$$

易知  $\mathbf{q}_i = 0$  的充要条件是  $\mathbf{r}_0, A\mathbf{r}_0, \dots, A^i\mathbf{r}_0$  线性相关。由  $A$ -正交化的要求, 有

$$\begin{aligned} \sigma_j^{(i-1)} &= -[A^i\mathbf{r}_0, \mathbf{q}_j] / [\mathbf{q}_j, \mathbf{q}_j] \\ &= -(\mathbf{q}_j, A^i\mathbf{r}_0) / (\mathbf{q}_j, \mathbf{q}_j), \end{aligned}$$

再由  $[\mathbf{q}_i, \mathbf{q}_j] = 0$ , 知负的  $\sigma_0^{(i-1)}\mathbf{q}_0 + \dots + \sigma_{i-1}^{(i-1)}\mathbf{q}_{i-1}$  是  $A^i\mathbf{r}_0$  在子空间  $\{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_{i-1}\}$  上的最佳逼近。

从正交化的过程可知

$$\{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_l\} \equiv \{\mathbf{r}_0, A\mathbf{r}_0, \dots, A^l\mathbf{r}_0\} = m_l.$$

对于向量  $\mathbf{x}_0 - \tilde{\mathbf{x}}$  在子空间  $m_l$  上找一个向量  $\mathbf{y}$ , 使得

$[x_0 - \tilde{x} + y, x_0 - \tilde{x} + y]$  达到极小. 利用定理 1.1 知道, 这样的  $y$  必须且只须满足方程

$$[x_0 - \tilde{x} + y, q_i] = 0, \quad i=0, 1, 2, \dots, l. \quad (17)$$

若记 
$$y = \sum_{i=0}^l \alpha_i q_i,$$

由 (17) 和  $q_0, q_1, \dots, q_l$  的正交性, 得到

$$\alpha_i = -[x_0 - \tilde{x}, q_i] / [q_i, q_i] = -(\mathbf{r}_0, q_i) / (Aq_i, q_i).$$

因为  $[x_0 - \tilde{x} + y, x_0 - \tilde{x} + y] = H'_1(x_0 + y),$

故知这样得到的  $y$  使得  $x_0 + y$  即为  $l+1$  步斜量法中得到的  $x_1$ , 从而有

$$\left| \tilde{x} - x_0 - \sum_{i=0}^l \alpha_i q_i \right| \leq E_{l+1} |\tilde{x} - x_0|.$$

因为  $l$  可以是任意非负整数, 注意到  $\alpha_i$  的计算公式不随  $l$  变化而变化, 因此我们可以得到如下的计算程式:

$$\begin{aligned} x_1 &= x_0 + \alpha_0 q_0, \\ \alpha_0 &= -(\mathbf{r}_0, q_0) / (Aq_0, q_0), \\ x_2 &= x_0 + \alpha_0 q_0 + \alpha_1 q_1 = x_1 + \alpha_1 q_1, \\ \alpha_1 &= -(\mathbf{r}_0, q_1) / (Aq_1, q_1), \\ x_3 &= x_0 + \alpha_0 q_0 + \alpha_1 q_1 + \alpha_2 q_2 = x_2 + \alpha_2 q_2, \\ \alpha_2 &= -(\mathbf{r}_0, q_2) / (Aq_2, q_2), \\ &\dots\dots\dots \\ x_k &= x_{k-1} + \alpha_{k-1} q_{k-1}, \\ \alpha_{k-1} &= -(\mathbf{r}_0, q_{k-1}) / (Aq_{k-1}, q_{k-1}). \end{aligned}$$

这一计算  $\{x_k\}$  的方法, 就是共轭斜量法, 因为用到了  $\mathbf{r}_0, A\mathbf{r}_0, \dots, A^l \mathbf{r}_0$  的  $A$ -正交化组  $q_0, q_1, \dots, q_l$ , 而两个向量  $A$ -正交又称为共轭, 所以称为共轭斜量法. 并且得到估计式:

$$\|x_k - \tilde{x}\| \leq \sqrt{\frac{\lambda_n}{\lambda_1}} E_k \|x_0 - \tilde{x}\|.$$

由正交化的过程中和计算  $w_k$  的过程中可以看到如果  $q_{k-1}=0$ , 那么计算就要中断, 但是恰巧在这时,  $x_{k-1}=\tilde{x}$ , 意味着计算已经完成. 下面来证明这一点, 首先由  $q_{k-1}=0$ , 知道  $r_0, Ar_0, \dots, A^{k-1}r_0$  是线性相关的, 因此  $x_0-\tilde{x}$  落在空间  $\{r_0, Ar_0, \dots, A^{k-2}r_0\}$  中, 也即  $x_0-\tilde{x}$  落在空间  $\{q_0, q_1, \dots, q_{k-2}\}$  中, 实际上从  $r_0, Ar_0, \dots, A^{k-1}r_0$  线性相关, 知存在常数  $\omega_l$ , 使  $\sum_{l=0}^{k-1} \omega_l A^l r_0 = 0$ , 等式两边作用  $A^{-1}$ , 即得

$$\sum_{l=0}^{k-1} \omega_l A^l (x_0 - \tilde{x}) = 0.$$

如果  $\omega_0 \neq 0$ , 就知  $x_0 - \tilde{x}$  可用  $r_0, Ar_0, \dots, A^{k-2}r_0$  来表示. 如果  $\omega_0 = 0, \omega_1 \neq 0$ , 可以再作用  $A^{-1}$ , 就得  $x_0 - \tilde{x}$  可用  $r_0, Ar_0, \dots, A^{k-3}r_0$  线性组合表示. 以此类推, 可知这一结论成立. 这样就知

$$\tilde{x} = x_0 + \beta_0 q_0 + \beta_1 q_1 + \dots + \beta_{k-2} q_{k-2},$$

但  $x_{k-1} = x_0 + \alpha_0 q_0 + \alpha_1 q_1 + \dots + \alpha_{k-2} q_{k-2},$

$$\begin{aligned} \text{而且 } 0 &\leq \left[ x_0 + \sum_{i=0}^{k-2} \alpha_i q_i - \tilde{x}, x_0 + \sum_{i=0}^{k-2} \alpha_i q_i - \tilde{x} \right] \\ &\leq \left[ x_0 + \sum_{i=0}^{k-2} \beta_i q_i - \tilde{x}, x_0 + \sum_{i=0}^{k-2} \beta_i q_i - \tilde{x} \right] = 0, \end{aligned}$$

$$\text{故必须 } x_{k-1} = x_0 + \sum_{i=0}^{k-2} \alpha_i q_i = x_0 + \sum_{i=0}^{k-2} \beta_i q_i = \tilde{x}.$$

这里还要指出一点: 从 (16) 那样的方法获得  $A$ -正交化组, 无论从计算量角度还是从存贮量角度来说都是花费太大, 因为在计算  $q_i$  时, 要用到  $q_0, q_1, \dots, q_{i-1}$ , 后面这  $i$  个向量都要存放在快速存贮器中, 而且不能指望这些向量是稀疏的, 因此得有  $n \times i$  个存贮单元, 当  $i$  较大时, 例如  $i = \frac{n}{2}$ , 那么就需  $\frac{1}{2} n^2$  个存贮单元, 这在实际计算中是太大的花费. 另外因为

用到  $q_0, q_1, \dots, q_{i-1}$ , 相应的要计算  $\sigma_0^{(i-1)}, \sigma_1^{(i-1)}, \dots, \sigma_{i-1}^{(i-1)}$ , 也要付出相应的计算量.

因此要设法, 在计算  $q_i$  时只用到少数几个向量. 为此利用  $\{q_0, q_1, \dots, q_{i-1}, Aq_{i-1}\} = \{r_0, Ar_0, \dots, A^i r_0\}$ , 于是可以将  $q_i$  表示成

$$q_i = Aq_{i-1} + \lambda_0^{(i-1)} q_0 + \lambda_1^{(i-1)} q_1 + \dots + \lambda_{i-1}^{(i-1)} q_{i-1},$$

利用  $[q_i, q_j] = 0, j=0, 1, 2, \dots, i-1$ , 知道

$$\begin{aligned}\lambda_j^{(i-1)} &= -[Aq_{i-1}, q_j] / [q_j, q_j] \\ &= -[q_{i-1}, Aq_j] / [q_j, q_j],\end{aligned}$$

$$\text{但 } Aq_j = q_{j+1} - \sum_{l=0}^j \lambda_l^{(j)} q_l,$$

$$\lambda_j^{(i-1)} = -\left[ q_{i-1}, q_{j+1} - \sum_{l=0}^j \lambda_l^{(j)} q_l \right] / [q_j, q_j],$$

可知当  $j+1 < i-1$  或  $j < i-2$  时

$$\lambda_j^{(i-1)} = 0,$$

$$\text{由此 } q_i = Aq_{i-1} + \lambda_{i-1}^{(i-1)} q_{i-1} + \lambda_{i-2}^{(i-1)} q_{i-2},$$

$$\text{并且 } \lambda_{i-1}^{(i-1)} = -[Aq_{i-1}, q_{i-1}] / [q_{i-1}, q_{i-1}],$$

$$\lambda_{i-2}^{(i-1)} = -[q_{i-1}, q_{i-1}] / [q_{i-2}, q_{i-2}],$$

这样每次计算  $q_i$  只要用到  $q_{i-1}, q_{i-2}$  两个向量, 节约了存储量和计算量.

新的计算程式如下:

取一个初始向量  $x_0$ , 计算  $r_0 = Ax_0 - b$ , 取  $q_0 = r_0$ ,

$$\begin{cases} \alpha_k = -(r_k, q_k) / (Aq_k, q_k), \\ x_{k+1} = x_k + \alpha_k q_k, \quad r_{k+1} = r_k + \alpha_k Aq_k, \\ q_{k+1} = Aq_k + \lambda_k^{(k)} q_k + \lambda_{k-1}^{(k)} q_{k-1}, \\ \lambda_k^{(k)} = -(Aq_k, Aq_k) / (Aq_k, q_k), \\ \lambda_{k-1}^{(k)} = -(Aq_k, q_k) / (Aq_{k-1}, q_{k-1}), \quad \lambda_{-1}^{(0)} = 0, \\ k = 0, 1, 2, \dots \end{cases} \quad (18)$$

共轭斜量法产生的序列  $\{x_k\}$ , 有一个重要的关系式

$$(r_k, r_j) = 0, \quad j = 0, 1, \dots, k-1, \quad (19)$$

实际上

$$\begin{aligned} (r_1, r_0) &= (r_0 + \alpha_0 A q_0, r_0) = (r_0, r_0) + \alpha_0 (A q_0, r_0) \\ &= (r_0, r_0) - (r_0, r_0) = 0. \end{aligned}$$

由  $x_k$  的定义知

$$(x_k - \tilde{x}, q_j) = 0, \quad j = 0, 1, \dots, k-1,$$

即  $(r_k, q_j) = 0, \quad j = 0, 1, \dots, k-1.$

只要证明  $r_j \in \{q_0, q_1, \dots, q_j\}$ , (19) 式就得到证明了。实际上

$$r_0 = q_0,$$

假设  $r_l \in \{q_0, q_1, \dots, q_l\}$  对于  $l = 0, 1, \dots, j-1$  都成立, 来证明  $r_j \in \{q_0, q_1, \dots, q_j\}$ . 因为

$$\begin{aligned} r_j &= r_{j-1} + \alpha_{j-1} A q_{j-1} \\ &= r_{j-1} + \alpha_{j-1} (q_j - \lambda_{j-1}^{(j-1)} q_{j-1} - \lambda_{j-2}^{(j-1)} q_{j-2}), \end{aligned}$$

因此  $r_j \in \{q_0, q_1, \dots, q_j\}$ , 从而证明了(19).

如果  $r_j = 0$ , 那么  $x_j = \tilde{x}$ ; 因为  $n$  维空间最多有  $n$  个非零彼此正交向量, 因此总有一个  $j \leq n$  使  $r_j = 0$ . 这也就是说共轭斜量法, 在理论上是有限步就可以结束的.

如果  $r_0, r_1, \dots, r_j$  全不为零, 则可以证明

$$\begin{aligned} \{r_0, r_1, \dots, r_j\} &= \{q_0, q_1, \dots, q_j\} \\ &= \{r_0, A r_0, \dots, A^j r_0\}, \end{aligned}$$

从而  $r_0, r_1, \dots, r_j$  是  $\{r_0, A r_0, \dots, A^j r_0\}$  的一组正交基. 前面已经证明了  $r_i \in \{q_0, q_1, \dots, q_j\}$ ,  $i = 0, 1, 2, \dots, j$ , 现在只要证明  $q_i \in \{r_0, r_1, \dots, r_j\}$ ,  $i = 0, 1, 2, \dots, j$ , 那么就有

$$\{r_0, r_1, \dots, r_j\} = \{q_0, q_1, \dots, q_j\}.$$

实际上从



$\mathbf{r}_i = \mathbf{r}_{i-1} + \alpha_{i-1}(\mathbf{q}_i - \lambda_{i-1}^{(i-1)} \mathbf{q}_{i-1} - \lambda_{i-2}^{(i-1)} \mathbf{q}_{i-2}), \quad i=1, \dots, j,$   
 知道  $\alpha_{i-1} \neq 0, i=1, \dots, j$ , 否则

$$\mathbf{r}_i = \mathbf{r}_{i-1}.$$

由  $(\mathbf{r}_i, \mathbf{r}_{i-1}) = 0$  得  $(\mathbf{r}_i, \mathbf{r}_i) = 0$ , 与  $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_j$  不全为零矛盾. 由  $\alpha_{i-1} \neq 0$ , 从而可知  $\mathbf{q}_i$  可由  $\mathbf{r}_i, \mathbf{r}_{i-1}, \dots, \mathbf{r}_0$  线性组合来表示, 即  $\mathbf{q}_i \in \{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_j\}, i=0, 1, 2, \dots, j$ .

现在我们可以利用  $\mathbf{r}_i$  的正交性, 来进一步减少  $A$ -正交化过程的计算量和存储量. 在(18)的计算过程中, 在计算得到  $\mathbf{r}_{k+1}$  后, 若  $\mathbf{r}_{k+1} \neq 0$  利用下式计算  $\mathbf{q}_{k+1}$ ,

$$\mathbf{q}_{k+1} = \mathbf{r}_{k+1} + \omega_0^{(k)} \mathbf{q}_0 + \omega_1^{(k)} \mathbf{q}_1 + \dots + \omega_k^{(k)} \mathbf{q}_k.$$

由  $[\mathbf{q}_{k+1}, \mathbf{q}_j] = 0, j=0, 1, \dots, k$ , 知

$$\omega_j^{(k)} = -[\mathbf{r}_{k+1}, \mathbf{q}_j] / [\mathbf{q}_j, \mathbf{q}_j],$$

$$[\mathbf{r}_{k+1}, \mathbf{q}_j] = (\mathbf{r}_{k+1}, A\mathbf{q}_j) = \left( \mathbf{r}_{k+1}, \frac{\mathbf{r}_{j+1} - \mathbf{r}_j}{\alpha_j} \right),$$

因此当  $j=0, 1, \dots, k-1$  时,  $\omega_j^{(k)} = 0$ , 这样

$$\mathbf{q}_{k+1} = \mathbf{r}_{k+1} + \omega_k^{(k)} \mathbf{q}_k,$$

$\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_k, k+1$  个向量中, 只有  $\mathbf{q}_k$  在计算  $\mathbf{q}_{k+1}$  时要用到, 又减少了一个向量. 此时

$$\begin{aligned} \omega_k^{(k)} &= -[\mathbf{r}_{k+1}, \mathbf{q}_k] / [\mathbf{q}_k, \mathbf{q}_k] \\ &= -(\mathbf{r}_{k+1}, \mathbf{r}_{k+1}) / \alpha_k [\mathbf{q}_k, \mathbf{q}_k], \end{aligned}$$

$$\begin{aligned} \alpha_k &= -(\mathbf{r}_k, \mathbf{q}_k) / (A\mathbf{q}_k, \mathbf{q}_k) \\ &= -(\mathbf{r}_k, \mathbf{r}_k + \omega_{k-1}^{(k-1)} \mathbf{q}_{k-1}) / (A\mathbf{q}_k, \mathbf{q}_k) \\ &= -(\mathbf{r}_k, \mathbf{r}_k) / (A\mathbf{q}_k, \mathbf{q}_k), \end{aligned}$$

于是

$$\omega_k^{(k)} = (\mathbf{r}_{k+1}, \mathbf{r}_{k+1}) / (\mathbf{r}_k, \mathbf{r}_k).$$

将  $\omega_k^{(k)}$  记为  $\lambda_k$ , 我们得到另外一个计算程式:

取一个初始向量  $x_0$ , 计算  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b}$ , 取  $\mathbf{q}_0 = \mathbf{r}_0$ ,

$$\begin{cases} \alpha_k = -(\mathbf{r}_k, \mathbf{r}_k) / (A\mathbf{q}_k, \mathbf{q}_k), \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{q}_k, \quad \mathbf{r}_{k+1} = \mathbf{r}_k + \alpha_k A\mathbf{q}_k, \\ \mathbf{q}_{k+1} = \mathbf{r}_{k+1} + \lambda_k \mathbf{q}_k, \\ \lambda_k = (\mathbf{r}_{k+1}, \mathbf{r}_{k+1}) / (\mathbf{r}_k, \mathbf{r}_k), \\ k = 0, 1, 2, \dots \end{cases} \quad (20)$$

(20)式也是流行的共轭斜量法的计算程式。

我们介绍了三种求  $A$ -正交化向量组的方法即(16), (18)和(20), 各种方法所求的  $\mathbf{q}_k$ , 长度  $\|\mathbf{q}_k\|$  可能会有差别, 但它们之间最多差一个常数因子。从计算量和存贮量节省来说(20)是最好了。使用(20)式计算  $\mathbf{x}_{k+1}$ , 需要用到  $\mathbf{x}_k, \mathbf{r}_k, \mathbf{q}_k, A\mathbf{q}_k$ , 一般它们都要存放在快速存贮器中, 这样要占用  $4n$  个存贮单元, 这一点比起其他求方程组的解的迭代法来说是比较费的。如果在每一步迭代中, 多计算一次  $A\mathbf{q}_k$ , 就可以不存放  $A\mathbf{q}_k$ , 即在计算

$$\alpha_k = -(\mathbf{r}_k, \mathbf{r}_k) / (A\mathbf{q}_k, \mathbf{q}_k)$$

中, 计算一次  $A\mathbf{q}_k$ , 作内积  $(A\mathbf{q}_k, \mathbf{q}_k)$  后, 不保留  $A\mathbf{q}_k$ , 在计算

$$\mathbf{r}_{k+1} = \mathbf{r}_k + \alpha_k A\mathbf{q}_k$$

时, 再计算一次  $A\mathbf{q}_k$ , 计算完  $\mathbf{r}_k$  后, 也不保留  $A\mathbf{q}_k$ , 这是一种增加计算量换得减少存贮量的措施。

共轭斜量法不但可用来求解系数矩阵为对称正定时的方程组, 而且还可以用来求解系数矩阵为对称非负定时的方程组, 也即系数矩阵  $A$  除了有正的特征值外, 还有零特征值时, 即行列式  $\det(A) = 0$  时, 也可使用共轭斜量法。考虑方程组

$$A\mathbf{x} = \mathbf{b} \quad (1)$$

的系数矩阵  $A$  的特征值  $\lambda_1, \lambda_2, \dots, \lambda_n$  中有  $\lambda_1, \lambda_2, \dots, \lambda_l$  是零而  $\lambda_{l+1} > 0$ ,  $\lambda_l$  对应的特征向量为  $\mathbf{y}_i$ 。于是(1)要有解的充要条件是

$$(\mathbf{b}, \mathbf{y}_i) = 0, \quad i=1, 2, \dots, l.$$

记集合

$$H_0 = \{\mathbf{x} \mid (\mathbf{x}, \mathbf{y}_i) = 0, i=1, 2, \dots, l\},$$

显然  $H_0$  是一个线性子空间, 并且  $\mathbf{b} \in H_0$ . 在  $H_0$  上引进函数  $(A\mathbf{x}, \mathbf{y}), \mathbf{x}, \mathbf{y} \in H_0$ , 容易验证  $(A\mathbf{x}, \mathbf{y})$  在  $H_0$  上满足内积的所有条件, 因此可以看作空间  $H_0$  上的内积.

因为(1)的解非唯一, 可以差一个零特征向量的线性组合  $\sum_{i=1}^l \beta_i \mathbf{y}_i$ . 因此对于任意取定的初始向量  $\mathbf{x}_0$ , 有一个解  $\tilde{\mathbf{x}}$ , 使得  $\mathbf{x}_0 - \tilde{\mathbf{x}} \in H_0$ , 自然有  $\mathbf{r}_0 = A\mathbf{x}_0 - \mathbf{b} \in H_0, A\mathbf{r}_0 \in H_0, \dots, A^k \mathbf{r}_0 \in H_0$ , 将在  $H_0$  上的内积  $(A\mathbf{x}, \mathbf{y})$  也记为  $[\mathbf{x}, \mathbf{y}]$ , 同样可以考虑多步斜量法, 使

$$\left[ \mathbf{x}_0 + \sum_{j=0}^{k-1} \beta_j A^j \mathbf{r}_0 - \tilde{\mathbf{x}}, \mathbf{x}_0 + \sum_{j=0}^{k-1} \beta_j A^j \mathbf{r}_0 - \tilde{\mathbf{x}} \right]$$

达到极小, 同样可以将  $\mathbf{r}_0, A\mathbf{r}_0, \dots, A^k \mathbf{r}_0$   $A$ -正交化, 就得到共轭斜量法, 计算程式与(20)一样, 并且有收敛速度的估计

$$\|\tilde{\mathbf{x}} - \mathbf{x}_k\| \leq \sqrt{\frac{\lambda_n}{\lambda_{l+1}}} \tilde{E}_k \|\tilde{\mathbf{x}} - \mathbf{x}_0\|,$$

其中 
$$\tilde{E}_k = \left[ T_k \left( \frac{\lambda_n + \lambda_{l+1}}{\lambda_n - \lambda_{l+1}} \right) \right]^{-1}.$$

在[5]中介绍了在实际问题中使用共轭斜量法解系数矩阵非负定时的方程组的一个例子.

当系数矩阵对称不定的情况, 即  $A$  有正的特征值, 也有负的特征值时, 一般不能再用(20)求解了. 对于这样的方程组, 用什么办法求解, 也是目前大家较关心的问题, 在第2章 §5 将介绍一种 Lanczos 算法. 当然在那种情况的各种算法的理论, 也没有象对称正定时共轭斜量法那样完全, 因此也是进一步需要研究的问题.

## §4 不完全分解、预处理共轭斜量法

上一节介绍的共轭斜量法，与其它求解方程组的方法比较起来，有下列几个很重要的优点：

1. 每步迭代需要计算的量，用到系数矩阵  $A$  的，只有  $A$  与向量的乘法  $Aq$ ，因此可以充分利用  $A$  的稀疏性。

2. 不要预先估计别的参数就可以计算，这一点不象契贝谢夫半迭代法、超松弛法等方法。

3. 每次迭代所需要计算，都是向量之间的运算，可充分为第四代计算机（向量运算计算机）所利用。

但是实际计算中，使用共轭斜量法，也会碰到问题，主要是收敛慢的问题，一方面是舍入误差的影响，更主要的是系数矩阵  $A$  的条件数太大。在 §1 中我们已经讲过条件数  $p = \frac{\lambda_n}{\lambda_1}$ ，

$$\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} = \frac{p+1}{p-1} = \frac{1 + \frac{1}{p}}{1 - \frac{1}{p}}, \quad \text{而共轭斜量法的收敛速度快慢依}$$

赖于

$$\left[ T_k \left( \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right) \right]^{-1} = \left[ T_k \left( \frac{1 + \frac{1}{p}}{1 - \frac{1}{p}} \right) \right]^{-1}$$

趋于 0 的速度。根据契贝谢夫多项式的定义，可知

$$\begin{aligned} T_k \left( \frac{1+\varepsilon}{1-\varepsilon} \right) &= \frac{1}{2} \left\{ \left( \frac{1+\sqrt{\varepsilon}}{1-\sqrt{\varepsilon}} \right)^k + \left( \frac{1+\sqrt{\varepsilon}}{1-\sqrt{\varepsilon}} \right)^{-k} \right\} \\ &\approx \frac{1}{2} \exp(2k\sqrt{\varepsilon}), \quad \text{当 } k\sqrt{\varepsilon} > 1, \varepsilon \rightarrow 0. \end{aligned}$$

当  $p$  很大时,  $1/p$ 、 $\sqrt{1/p}$  都很小, 因此  $\exp(-2k\sqrt{1/p})$  收敛于 0 很慢.

很早以前就有人提出, 对于方程组

$$Ax = b \quad (1)$$

作等价变换, 变成

$$Bx = f, \quad (21)$$

或者变成

$$By = g, \quad (22)$$

只要  $x$  能够容易从  $y$  获得, 而  $B$  仍然保持对称正定, 并且  $B$  的条件数  $p(B)$  要比  $A$  的条件数  $p(A)$  小, 这样解方程 (21) 或 (22) 要比解方程 (1) 有利. 当然使用共轭斜量法于 (21)、(22) 收敛速度也会比对于 (1) 的要快. 这就是预处理 (Precondition) 的思想. 但是很多预处理, 尽管条件数改善了, 但乃不象  $A$  那样具有稀疏性, 并且增加很多预处理的计算量, 这可参见 [6].

1977 年 J. A. Meijerink 和 A. Van der Vorst [7] 提出一种称为不完全分解 (Incomplete decomposition) 的办法, 将  $A$  分解成

$$A = LL^T + R,$$

这里  $L$  是下三角阵; 一方面使  $LL^T$  尽可能接近  $A$ , 另一方面使  $L$  保持跟  $A$  一样的稀疏性, 或者具有其他形状的稀疏性.

完全分解是对矩阵  $A$  进行三角分解  $A = LL^T$ , 不完全分解是对矩阵  $A - R$  进行三角分解  $LL^T$ . 对于不完全分解, 因为有矩阵  $R$ , 可以变化, 因此  $L$  中那些元素为 0, 可以预先规定, 不过这种规定也不能完全任意, 还须考虑到  $LL^T$  接近  $A$ , 但是这是不易检验的, 因此实际计算时, 就考虑使  $R$  有较多零元素.

为了明白具体做法, 举  $4 \times 4$  对称正定矩阵  $A$  的不完全分解为例:  $A = (a_{ij})$ ,  $R = (r_{ij})$ , 要求分解成  $LL^T$ , 其中  $L$  的形状如下:

$$L = \begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ 0 & l_{32} & l_{33} & \\ 0 & 0 & l_{43} & l_{44} \end{pmatrix},$$

即  $l_{31} = l_{41} = l_{42} = 0$ , 比较

$$LL^T = (A - R),$$

有  $l_{11}^2 = a_{11} - r_{11}$ , 因为  $A$  正定, 故  $a_{11} > 0$ , 因此可取  $r_{11} = 0$ , 于是  $l_{11} = \sqrt{a_{11}}$ ; 再由  $l_{11}l_{21} = (a_{12} - r_{12})$ , 可以取  $r_{12} = 0$ ,  $l_{21} = a_{12}/l_{11}$ . 由  $l_{11}l_{31} = (a_{31} - r_{31})$ , 因为  $l_{31} = 0$ , 因此必须取  $r_{31} = a_{31}$ ; 同样由  $l_{41} = 0$ , 从  $l_{11}l_{41} = (a_{41} - r_{41})$  确定  $r_{41} = a_{41}$ . 然后由

$$l_{21}^2 + l_{22}^2 = a_{22} - r_{22},$$

$$l_{22}^2 = a_{22} - l_{21}^2 - r_{22},$$

如果  $a_{22} - l_{21}^2 > 0$ , 那么可取  $r_{22} = 0$ , 在完全分解时, 当  $A$  对称正定必有  $a_{22} - l_{21}^2 > 0$ , 对于不完全分解, 就不一定能保证  $a_{22} - l_{21}^2 > 0$ , 如果  $a_{22} - l_{21}^2 \leq 0$ , 取一个绝对值小的负数  $r_{22}$ , 使  $a_{22} - l_{21}^2 - r_{22} > 0$ , 这样

$$l_{22} = \sqrt{a_{22} - l_{21}^2 - r_{22}},$$

由  $l_{21}l_{31} + l_{22}l_{32} = (a_{32} - r_{32})$ , 可取  $r_{32} = 0$ ,  $l_{32} = a_{32}/l_{22}$ , 由  $l_{31}l_{41} + l_{32}l_{42} = a_{42} - r_{42}$ , 要  $l_{42} = 0$ , 必须  $r_{42} = a_{42}$ , 由  $l_{31}^2 + l_{32}^2 + l_{33}^2 = a_{33} - r_{33}$ ,  $l_{33}^2 = a_{33} - l_{31}^2 - r_{33}$ , 同样如果  $a_{33} - l_{31}^2 > 0$ , 那么  $r_{33} = 0$ , 否则取  $r_{33}$  为一个使  $a_{33} - l_{31}^2 - r_{33} > 0$  的负数, 于是  $l_{33} = \sqrt{a_{33} - l_{31}^2 - r_{33}}$ , 由  $l_{31}l_{41} + l_{32}l_{42} + l_{33}l_{43} = a_{43} - r_{43}$ , 取  $r_{43} = 0$ ,  $l_{43} = a_{43}/l_{33}$ ; 最后由

$$l_{44}^2 = a_{44} - l_{43}^2 - r_{44}$$

决定  $r_{44}$  和  $l_{44} = \sqrt{a_{44} - l_{43}^2 - r_{44}}$ . 按此得到的

$$R = \begin{pmatrix} 0 & 0 & a_{31} & a_{41} \\ 0 & r_{22} & 0 & a_{42} \\ a_{31} & 0 & r_{33} & 0 \\ a_{41} & a_{42} & 0 & r_{44} \end{pmatrix},$$

如果  $A$  中的元素  $a_{31}, a_{41}, a_{42}$  有的是 0, 那么  $R$  中就会有更多的零. 由此可以设想,  $A$  稀疏时,  $R$  会有很多零, 对于对角元素  $r_{22}, r_{33}, r_{44}$ , 一般不能保证为 0, 但是文献 [7] 中证明, 当  $A$  是对称正定的  $M$ -阵时, 即

$a_{ij} \leq 0 (i \neq j), a_{ii} > 0, A^{-1}$  元素全为正时, 不管  $L$  的形状如何, 所有  $r_{ii} = 0$ .

对于  $A$  不是  $M$ -阵时, 不能保证所有  $r_{ii} = 0$ , 例如

$$A = \begin{pmatrix} 3 & -2 & 0 & 2 \\ -2 & 3 & -2 & 0 \\ 0 & -2 & 3 & -2 \\ 2 & 0 & -2 & 3 \end{pmatrix},$$

要求  $L$  与  $A$  有相同的稀疏性, 即  $l_{31} = l_{42} = 0$ , 此时

$$L = \begin{pmatrix} \sqrt{3} & & & 0 \\ -2/\sqrt{3} & \sqrt{5}/3 & & \\ 0 & -2\sqrt{3}/\sqrt{5} & \sqrt{3}/\sqrt{5} & \\ 2/\sqrt{3} & 0 & -2\sqrt{5}/\sqrt{3} & a \end{pmatrix},$$

$$R = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 4/3 \\ 0 & 0 & 0 & 0 \\ 0 & 4/3 & 0 & -(5+a^2) \end{pmatrix}.$$

所谓不完全分解预处理共轭斜量法, 即是先作  $A$  的某种不完全分解

$$A = LL^T + R,$$

再将方程组

$$Ax = b$$

化成

$$L^{-1}AL^{-T}L^Tx = L^{-1}b,$$

然后对方程组

$$L^{-1}AL^{-T}y = L^{-1}b \quad (23)$$

使用共轭斜量法, 求得  $y_k$ , 再由  $L^Tx = y_k$  求得  $x_k$ . 因为  $A$  对称正定, 易知  $L^{-1}AL^{-T}$  也是对称正定, 因此对 (23) 使用共轭斜量法, 是可行的.

具体计算的程式如下:

先对方程组 (1) 的  $A$  进行不完全分解, 得  $L$ .

取  $x_0$ , 算  $r_0 = Ax_0 - b$ , 求  $p_0 = (LL^T)^{-1}r_0$ ,

$$\alpha_k = -((LL^T)^{-1}r_k, r_k) / (Ap_k, p_k),$$

$$x_{k+1} = x_k + \alpha_k p_k,$$

$$r_{k+1} = r_k + \alpha_k Ap_k,$$

$$\lambda_k = ((LL^T)^{-1}r_{k+1}, r_{k+1}) / ((LL^T)^{-1}r_k, r_k),$$

$$p_{k+1} = (LL^T)^{-1}r_{k+1} + \lambda_k p_k,$$

$$k = 0, 1, 2, \dots,$$

以上就是不完全分解、预处理共轭斜量法的计算程式, 这个方法简记为 IOOG 法.

程式中出现  $(LL^T)^{-1}r_k$ , 表示要解一个方程组

$$(LL^T)y = r_k,$$

但因为  $L$  是下三角阵, 又是稀疏的, 因此计算量和存贮量都可以不大.

有很多实际计算的报告, 指出这一方法非常成功, 可参见 [7]、[8], 例如 [8] 中指出对于一个激光聚变 (Laser fusion) 方



程

$$\frac{\partial f}{\partial t} = \nabla D \nabla f + r(f^e - f),$$

这里  $D$ 、 $r$ 、 $f^e$  是给定的非负函数。边界条件是

$$\mathbf{n} \cdot \nabla f = 0 \quad \text{和} \quad f = 0$$

两部分的混合。

用差分方法解此方程, 用隐式 5 点差分格式, 共取 555 个节点, 化成一个线性方程组

$$A\mathbf{x} = \mathbf{b},$$

系数矩阵  $A$  的条件数  $p = 2 \times 10^{10}$ 。

如果记  $\tilde{\mathbf{x}}$  为正确解, 要求计算到

$$\varepsilon = \|\tilde{\mathbf{x}} - \mathbf{x}\| / \|\tilde{\mathbf{x}}\| = 10^{-6},$$

用下列五种方法进行计算, 所要求的迭代次数比较如下:

1. ICCG 方法

$$A = LL^T + R,$$

其中  $L$  的稀疏性跟  $A$  同, 记为 ICCG(0), 需 25 次迭代。

2. 取最佳松弛因子的块超松弛法, 需 765 次迭代。

3. 隐式交替方向法 1, 需 10200 次迭代。

4. 隐式交替方向法 2, 需 4750 次迭代。

5. Gauss-Seidel 法, 需 208000 次迭代。

对于这个问题, ICCG(0) 比 Gauss-Seidel 法要快了 8000 倍。

这个例子说明 ICCG 法在实用中有很惊人的成绩。但是迄今对这一方法的理论分析, 还没有重要结果。究竟  $L^{-1}AL^{-T}$  的条件数比  $A$  的条件数, 一般情况下小多少? 当  $L$  取怎样的稀疏性时较好? 当  $A$  不是  $M$ -阵时,  $l_u$  取什么数较好? 这些问题都没有解决。目前有很多学者, 正在研究 ICCG 方法。

## 对称三对角矩阵

对称三对角矩阵, 是一种运算, 存放都比较简单的矩阵. 一个对称矩阵可以通过有限步计算, 正交相似变换成一个对称三对角矩阵, 因此很多求对称矩阵特征值的方法, 第一步先把矩阵化成对称三对角阵, 然后再来求三对角阵的特征值. 有些求解系数矩阵对称的线性代数方程组的方法, 也是先把系数矩阵, 归化成对称三对角阵, 然后再来解系数矩阵为对称三对角阵的方程组. 因此在矩阵计算中, 对称三对角矩阵, 已成为一种有用的工具. 在这一章介绍对称三对角矩阵的性质, 一方面为了后面两章的应用, 另外方面也有它本身的独立意义. 另外还介绍了解对称线性方程组的 Lanczos 算法.

### §1 Jacobi 矩阵

$$\text{设 } A = \begin{pmatrix} b_1 & c_1 & & & 0 \\ a_1 & b_2 & c_2 & & \\ & a_2 & b_3 & c_3 & \\ & & \ddots & \ddots & \ddots \\ 0 & & & \ddots & c_{n-1} \\ & & & a_{n-1} & b_n \end{pmatrix}$$

称为三对角矩阵, 如果  $a_i$ 、 $b_i$ 、 $c_i$  都是实数, 且  $a_i c_i > 0$  则称  $A$  为 Jacobi 矩阵.

从它的形状可知, 它的  $k$  阶顺序主子阵也是三对角阵, 记为  $A_k$ , 显然  $A_1 = b_1$ ,  $A_n = A$ .

设  $A_k$  的特征多项式为  $\varphi_k(\lambda)$ , 则

$$\varphi_k(\lambda) = \det(\lambda I - A_k)$$

$$= \begin{vmatrix} \lambda - b_1 & -c_1 & & & 0 \\ -a_1 & \lambda - b_2 & -c_2 & & \\ & \ddots & \ddots & \ddots & \\ 0 & & -a_{k-1} & \lambda - b_k & \end{vmatrix}$$

$$= (\lambda - b_k) \varphi_{k-1}(\lambda)$$

$$+ c_{k-1} \begin{vmatrix} \lambda - b_1 & -c_1 & & & 0 \\ -a_1 & \lambda - b_2 & -c_2 & & \\ & \ddots & \ddots & \ddots & \\ 0 & & -a_{k-3} & \lambda - b_{k-2} & -c_{k-2} \\ & & & 0 & -a_{k-1} \end{vmatrix}$$

$$= (\lambda - b_k) \varphi_{k-1}(\lambda) - a_{k-1} c_{k-1} \varphi_{k-2}(\lambda).$$

令  $\varphi_0(\lambda) = 1$ , 所以对  $k=2, 3, \dots, n$  成立关系式

$$\varphi_k(\lambda) = (\lambda - b_k) \varphi_{k-1}(\lambda) - a_{k-1} c_{k-1} \varphi_{k-2}(\lambda). \quad (1)$$

**定理 2.1** Jacobi 阵的特征多项式序列

$$\varphi_n(\lambda), \varphi_{n-1}(\lambda), \dots, \varphi_1(\lambda), \varphi_0(\lambda) \quad (2)$$

是任何区间  $[a, b]$  内的 Sturm 序列 (见 [9]).

**证明**

1. 序列最后一个多项式为  $\varphi_0(\lambda) = 1$ , 没有根.

2. 序列中相邻两个多项式  $\varphi_k(\lambda)$ ,  $\varphi_{k-1}(\lambda)$  没有公根. 否则由递推公式 (1) 知, 这个根也是  $\varphi_{k-2}(\lambda)$  的根, 从而可以推

知它也是  $\varphi_0(\lambda)$  的根, 此与  $\varphi_0(\lambda)$  无根矛盾.

3. 若  $\lambda_0$  是  $\varphi_{k-1}(\lambda)$  的根, 则  $\varphi_k(\lambda_0)$  与  $\varphi_{k-2}(\lambda_0)$  反号. 这是因为  $\varphi_k(\lambda_0) = -c_{k-1}a_{k-1}\varphi_{k-2}(\lambda_0)$ , 而  $c_{k-1}a_{k-1} > 0$  故  $\varphi_k(\lambda_0)$  与  $\varphi_{k-2}(\lambda_0)$  反号.

序列(2)具有上述三个性质, 因此为任何区间  $[a, b]$  内的一个 Sturm 序列. 证毕.

现在来分析序列(2)的零点分布.

1.  $\varphi_0(\lambda) = 1$ , 说明  $\varphi_0(\lambda)$  在整个实数轴  $(-\infty, \infty)$  上都保持正号.

2.  $\varphi_1(\lambda) = \lambda - b_1$ , 有一个零点, 即  $x_1^{(1)} = b_1$ ,  $\varphi_1(\lambda)$  在  $\lambda < b_1$  时为负, 在  $\lambda > b_1$  时为正.

3.  $\varphi_2(\lambda) = (\lambda - b_2)\varphi_1(\lambda) - a_1c_1\varphi_0(\lambda)$  是一个二次多项式, 首项系数为正, 因此在  $\lambda \rightarrow \pm\infty$  时,  $\varphi_2(\lambda)$  为正, 但在  $\lambda = x_1^{(1)}$  处  $\varphi_2(x_1^{(1)}) = -a_1c_1$  为负, 这说明  $\varphi_2(\lambda)$  有二个零点  $x_1^{(2)}$ 、 $x_2^{(2)}$ , 并且

$$x_1^{(2)} < x_1^{(1)} < x_2^{(2)}.$$

4. 一般地说, 可以用数学归纳法证明如下性质:  $\varphi_{k-1}(\lambda)$  有  $k-1$  个实单根  $x_1^{(k-1)}$ 、 $x_2^{(k-1)}$ 、 $\dots$ 、 $x_{k-1}^{(k-1)}$ ,  $\varphi_k(\lambda)$  有  $k$  个实单根  $x_1^{(k)}$ 、 $x_2^{(k)}$ 、 $\dots$ 、 $x_k^{(k)}$ , 并成立下列关系

$$x_1^{(k)} < x_1^{(k-1)} < x_2^{(k)} < x_2^{(k-1)} < \dots < x_{k-1}^{(k)} < x_{k-1}^{(k-1)} < x_k^{(k)},$$

实际上, 假如这一性质对  $k$  成立, 我们证明这一性质对  $k+1$  也成立.

考察  $\varphi_{k-1}(\lambda)$  在  $(x_i^{(k)} - \varepsilon, x_i^{(k)} + \varepsilon)$  中的符号, 这里  $\varepsilon$  是充分小的正数.

当  $k$  是奇数时,  $k-1$  是偶数,  $\varphi_{k-1}(\lambda)$  在  $(-\infty, x_1^{(k-1)})$  中是正的,  $\varphi_{k-1}(\lambda)$  在  $(x_1^{(k-1)}, x_2^{(k-1)})$  中是负的, 一般可知  $\varphi_{k-1}(\lambda)$

在  $(x_{i-1}^{(k-1)}, x_i^{(k-1)})$  中的符号同  $(-1)^{i-1}$ . 因为  $(x_i^{(k)} - \varepsilon, x_i^{(k)} + \varepsilon) \subset (x_{i-1}^{(k-1)}, x_i^{(k-1)})$ , 因此当  $k$  是奇数时  $\varphi_{k-1}(\lambda)$  在  $(x_i^{(k)} - \varepsilon, x_i^{(k)} + \varepsilon)$  中的符号同  $(-1)^{i-1}$ .

当  $k$  是偶数时,  $\varphi_{k-1}(\lambda)$  在  $(-\infty, x_1^{(k-1)})$  中的符号为负, 在  $(x_1^{(k-1)}, x_2^{(k-1)})$  中符号为正, 一般, 在  $(x_{i-1}^{(k-1)}, x_i^{(k-1)})$  中的符号同  $(-1)^i$ , 因此  $\varphi_{k-1}(\lambda)$  在  $(x_i^{(k)} - \varepsilon, x_i^{(k)} + \varepsilon)$  中的符号同  $(-1)^i$ .

因为  $\varphi_{k+1}(x_i^{(k)})$  的符号与  $\varphi_{k-1}(x_i^{(k)})$  的相反, 因此当  $k$  是奇数时,  $\varphi_{k+1}(x_i^{(k)})$  的符号同  $(-1)^i$ , 这样  $\varphi_{k+1}(\lambda)$  在  $(x_i^{(k)}, x_{i+1}^{(k)})$  中至少有一个根, 于是在  $(x_1^{(k)}, x_k^{(k)})$  中至少有  $k-1$  个根, 因为  $k$  是奇数,  $\varphi_{k+1}(x_1^{(k)})$  是负的,  $\varphi_{k+1}(x_k^{(k)})$  也是负的, 而  $\varphi_{k+1}(\lambda)$  是偶数次多项式, 在  $\lambda \rightarrow \pm\infty$  处都是正的, 因此在  $(-\infty, x_1^{(k)})$  中至少有一个根, 在  $(x_k^{(k)}, \infty)$  中也至少有一个根,  $\varphi_{k+1}(\lambda)$  至多有  $k+1$  个根, 于是知道在  $(x_i^{(k)}, x_{i+1}^{(k)})$  中有且只有一个根. 记为  $x_{i+1}^{(k+1)}$ , 同样在  $(-\infty, x_1^{(k)})$  中有且只有一个根记为  $x_1^{(k+1)}$ , 在  $(x_k^{(k)}, \infty)$  中有且只有一个根记为  $x_{k+1}^{(k+1)}$ . 这就证明了当  $k$  是奇数时所述性质对  $k+1$  也成立.

对于  $k$  是偶数的情况类似的可以证明. 我们有

**定理 2.2** Jacobi 矩阵的特征多项式序列  $\varphi_0(\lambda), \varphi_1(\lambda), \dots, \varphi_n(\lambda)$  中任意一个多项式  $\varphi_k(\lambda)$ , 有  $k$  个实单根  $x_1^{(k)}, x_2^{(k)}, \dots, x_k^{(k)}$ , 并且它们与  $\varphi_{k-1}(\lambda)$  的  $k-1$  个实单根  $x_1^{(k-1)}, x_2^{(k-1)}, \dots, x_{k-1}^{(k-1)}$ , 有如下隔离性质

$$x_1^{(k)} < x_1^{(k-1)} < x_2^{(k)} < x_2^{(k-1)} < \dots < x_{k-1}^{(k-1)} < x_k^{(k)}.$$

## § 2 对称三对角矩阵的唯一归化定理

形如

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \ddots & \ddots & \ddots & \\ 0 & & & \beta_{n-1} & \\ & & & \beta_{n-1} & \alpha_n \end{pmatrix}$$

的方阵称为对称三对角矩阵, 当  $\beta_i \neq 0$  ( $i=1, 2, \dots, n-1$ ) 时,  $T$  称为不可约。显然它是特殊的 Jacobi 阵, 因此关于 Jacobi 阵的结论, 它都具有。如果有  $\beta_i = 0$ , 那么矩阵是可约的 [9, p. 230]。可约对称三对角阵的方程组求解问题和特征值问题, 都可化成几个低阶的不可约矩阵的问题来讨论。本章以后几节将进一步介绍对称三对角矩阵的重要性质。

首先指出, 任意一个实对称矩阵, 都可以通过正交相似变换, 变成对称三对角矩阵。

取  $A = (a_{ij})$  是任意一个实对称矩阵, 它可以通过 Householder 镜像变换  $Q_1$  使

$$Q_1 A Q_1^T = \begin{pmatrix} a_{11} & \tau & 0 & \cdots & 0 \\ \tau & & & & \\ 0 & & A_2 & & \\ \vdots & & & & \\ 0 & & & & \end{pmatrix},$$

其中  $Q_1 = I - 2w_1 w_1^T = Q_1^T$ ,

$$w_1^T = (0, a_{21} - \tau, a_{31}, \dots, a_{n1})/h,$$

$$\tau = -(\operatorname{sign} a_{21}) \sqrt{\sum_{j=2}^n a_{j1}^2}, \quad h = \sqrt{2\tau(\tau - a_{21})},$$

熟知  $Q_1$  是正交变换。  $\tau$  之所以要取成与  $-a_{21}$  同号, 目的是使  $h$  尽可能大一些, 有利于计算的稳定。 如果记

$$A_2 = \begin{pmatrix} a'_{22} & a'_{23} & \cdots & a'_{2n} \\ a'_{32} & a'_{33} & \cdots & a'_{3n} \\ \cdots & \cdots & \cdots & \cdots \\ a'_{n2} & a'_{n3} & \cdots & a'_{nn} \end{pmatrix},$$

则再通过镜像变换

$$Q_2 = I - 2w_2 w_2^T,$$

$w_2$  是由  $A_2$  的第一列构成

$$w_2^T = (0, 0, a'_{32} - \tau', a'_{42}, \cdots, a'_{n2})/h',$$

其中  $\tau' = -(\text{sign } a'_{32}) \sqrt{\sum_{j=2}^n (a'_{j2})^2}$ ,  $h' = \sqrt{2\tau'(\tau' - a'_{32})}$ .

$$Q_2 Q_1 A Q_1^T Q_2^T = \begin{pmatrix} a_{11} & \tau & 0 & 0 & \cdots & 0 \\ \tau & a'_{22} & \tau' & 0 & \cdots & 0 \\ 0 & \tau' & & & & \\ 0 & 0 & & & & \\ \vdots & \vdots & & & A_3 & \\ 0 & 0 & & & & \end{pmatrix},$$

$Q_2 Q_1$  仍是正交变换。依此类推, 经过  $Q_1, Q_2, \cdots, Q_{n-2}$  镜像变换, 可以将  $A$  相似变换成一个对称三对角矩阵。

也可以通过平面旋转变换, 将对称矩阵  $A$  正交相似变换成对称三对角矩阵。记

$$T_{l,k} = \begin{pmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & c & & s \\ & & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \\ & & & -s & & c & & \\ & & & & & & 1 & \\ & & & & & & & \ddots \\ & & & & & & & & 1 \end{pmatrix} \begin{matrix} l \text{ 行} \\ \\ \\ k \text{ 行} \end{matrix}$$

$l \text{ 列} \quad \quad k \text{ 列}$

$s^2 + c^2 = 1$ , 称为  $(l, k)$  平面上的旋转矩阵, 它也是一个正交阵.

$$B = (b_{ij}) = T_{lk} A T_{lk}^T,$$

称为将  $A$  通过  $(l, k)$  平面旋转相似变换成  $B$ , 此时  $B$  的元素与  $A$  的元素之间的关系是

$$b_{li} = b_{il} = a_{li}c + a_{ik}s,$$

$$b_{ki} = b_{ik} = -a_{li}s + a_{ik}c, \quad i \neq l, k;$$

$$b_{ll} = a_{ll}c^2 + 2a_{lk}sc + a_{kk}s^2,$$

$$b_{kk} = a_{ll}s^2 - 2a_{lk}sc + a_{kk}c^2,$$

$$b_{lk} = b_{kl} = (a_{kk} - a_{ll})sc + a_{lk}(c^2 - s^2),$$

其余的  $b_{ij} = a_{ij}$  (当  $i \neq l, k$ , 且  $j \neq l, k$ ).

于是可知  $B$  仍是对称, 跟  $A$  相比只改变二行二列元素. 如果  $i \neq l, k$ ,  $a_{li} = a_{ik} = 0$ , 那么通过  $(l, k)$  平面的旋转后, 仍有  $b_{li} = b_{ik} = 0$ ; 如果  $a_{li}^2 + a_{ik}^2 \neq 0$ , 则取  $s = a_{ik} / \sqrt{a_{li}^2 + a_{ik}^2}$ ,  $c = a_{li} / \sqrt{a_{li}^2 + a_{ik}^2}$ , 可使  $b_{ik} = 0$ .

利用上述性质, 提出两种方案将一个实对称矩阵, 正交相似变换成一个对称三对角矩阵.

1. 当矩阵是满秩矩阵时. 将下列足标的元素, 依次逐个化成零,

$$\begin{aligned} &(1, n), (1, n-1), \dots, (1, 4), (1, 3), \\ &(2, n), (2, n-1), \dots, (2, 4), \\ &\dots\dots\dots \\ &(n-2, n), \end{aligned} \quad (3)$$

而在化足标如  $(i, j)$  的元素为 0 时, 采用的旋转矩阵为  $T_{j-1, j}$ , 即在  $(j-1, j)$  平面上作旋转, 如果当时的矩阵为  $\tilde{A} = (\tilde{a}_{pq})$ , 则取

$$s = \tilde{a}_{ij} / \sqrt{\tilde{a}_{i, j-1}^2 + \tilde{a}_{ij}^2}, \quad c = \tilde{a}_{i, j-1} / \sqrt{\tilde{a}_{i, j-1}^2 + \tilde{a}_{ij}^2},$$

如果  $\tilde{a}_{ij} = 0$ , 那么  $T_{j-1, j} = I$ , 此时这个旋转可以不做. 这样至



多进行了  $\frac{(n-1)(n-2)}{2}$  次旋转相似变换后, 就将  $A$  化成一个

对称三对角矩阵. 即

$$T_{n-1,n} \cdots T_{n-2,n-1} T_{n-1,n} A T_{n-1,n}^T T_{n-2,n-1}^T \cdots T_{n-1,n}^T = H$$

是一个对称三对角矩阵, 上述表示式中, 同样记号  $T_{j-1,j}$  会出现不止一次, 例如  $T_{n-1,n}$  出现  $n-2$  次, 但每次出现的  $T_{n-1,n}$ , 它的  $s$ 、 $c$  可能是不同的, 其它的  $T_{ij}$  也是这样.

这个过程比起镜像变换的办法来, 计算量可能要大一倍. 因此对于满秩矩阵的情况, 一般使用镜像变换的办法.

2. 当矩阵是带形时. 矩阵元素  $a_{ij}$ , 当  $j > i$  时称它在第  $j-i$  条超对角线 (Superdiagonal) 上,  $a_{1n}$  在第  $n-1$  条超对角线上,  $a_{12}$  在第 1 条超对角线上. 所谓带形对称矩阵就是存在一个自然数  $m < n$ , 对任何  $k > m$ , 第  $k$  条超对角线上的元素全为 0. 当然, 有兴趣的是  $m \ll n$  的情况. 对于这样的  $m$ , 称矩阵为带宽是  $2m+1$  的带形矩阵.

对于这种带形矩阵, 通过平面旋转, 逐个把所有三对角线外的元素化成零. 化的次序是先对第  $m$  条超对角线上的所有元素, 将它们全部化成零后, 再对第  $m-1$  条超对角线上的所有元素, 最后对第 2 条超对角线上的元素. 在化每条超对角线上元素为零时, 所取元素的次序, 是按照行的次序; 例如在化第  $k$  条超对角线上元素时, 所取元素的足标次序为  $(1, k+1), (2, k+2), \dots, (n-k, n)$ .

在消去足标为  $(i, j)$  的元素时, 使用的平面旋转为  $T_{j-1,j}$ , 设此时的矩阵为  $\tilde{A} = (\tilde{a}_{pq})$ , 则取

$$s = \tilde{a}_{ij} / \sqrt{\tilde{a}_{ij-1}^2 + \tilde{a}_{ij}^2},$$

$$c = \tilde{a}_{ij-1} / \sqrt{\tilde{a}_{ij-1}^2 + \tilde{a}_{ij}^2},$$

如果  $\tilde{a}_{ij} = 0$ , 此时  $T_{j-1,j} = I$ , 不必进行这一次旋转. 相似变换

后的矩阵

$$B = T_{j-1,j} \tilde{A} T_{j-1,j}^T = (b_{pq}),$$

就有  $b_{ij} = 0$ .

不过必须特别指出, 在消去元素  $\tilde{a}_{ij}$  时, 次序大于  $j-i$  的超对角线上元素, 已经化成零了. 即  $\tilde{A}$  中的元素  $\tilde{a}_{l,k}$  当  $k-l > j-i$  时,  $\tilde{a}_{l,k} = 0$ , 但是对于元素  $\tilde{a}_{j,2j-i}$ , 如果  $2j-i \leq n$ , 矩阵中是有这个元素的, 它是在第  $j-i$  条超对角线上的, 因为它的行足标  $j > i$ , 因此它可能不为 0. 这样  $B$  的元素

$$b_{j-1,2j-i} = \tilde{a}_{j-1,2j-i}c + \tilde{a}_{j,2j-i}s = \tilde{a}_{j,2j-i}s \neq 0.$$

$\tilde{a}_{j-1,2j-i}$  是在第  $j-i+1$  条超对角线上, 所以为零. 于是在第  $j-i+1$  条超对角线上产生了新的非零元  $b_{j-1,2j-i}$ . 尽管产生了新的非零元, 但是要注意, 这个新的非零元的列足标比被消的对象  $\tilde{a}_{ij}$  的列足标要大  $j-i$ . 这一点使人们想到再用  $T_{2j-i-1,2j-i}$  将  $b_{j-1,2j-i}$  化成零. 即得到

$$A = T_{2j-i-1,2j-i} B T_{2j-i-1,2j-i}^T = (\tilde{a}_{pq}),$$

$$\begin{aligned} \text{此时取} \quad s &= b_{j-1,2j-i} / \sqrt{b_{j-1,2j-i}^2 + b_{j-1,2j-i-1}^2}, \\ c &= b_{j-1,2j-i-1} / \sqrt{b_{j-1,2j-i}^2 + b_{j-1,2j-i-1}^2}, \end{aligned}$$

就有  $\tilde{a}_{j-1,2j-i} = 0$ .

同样如果  $2j-i + (j-i) = 3j-2i \leq n$ , 而  $b_{2j-i,3j-2i}$  又不为零, 则又产生一个新的非零元  $\tilde{a}_{2j-i-1,3j-2i}$ , 不过列的足标又增加  $j-i$ ; 这样产生一个新的非零元, 用平面旋转将它化成零, 新的非零元又产生, 不过列足标增大  $j-i$ , 因为列足标最大是  $n$ . 因此当  $\left[ \frac{n-j}{j-i} \right]$  步后, 就不会再有新的非零元产生. 这样的办法, 消去一个  $\tilde{a}_{ij}$ , 总共要化  $\left[ \frac{n-j}{j-i} \right] + 1$  次平面旋转变换.

但这样的过程, 最大的优点是保持带宽, 矩阵的带宽以外部分不必占用存贮单元. 而镜象变换, 一般就不能保持带宽, 在这

一点上比镜像变换优越。很多实际问题产生的矩阵，都是稀疏带形矩阵，因此这个办法是很有实用价值的。这个方法最早由 H. R. Schwarz 给出，见[24]。图 2 给出这个过程的示意。

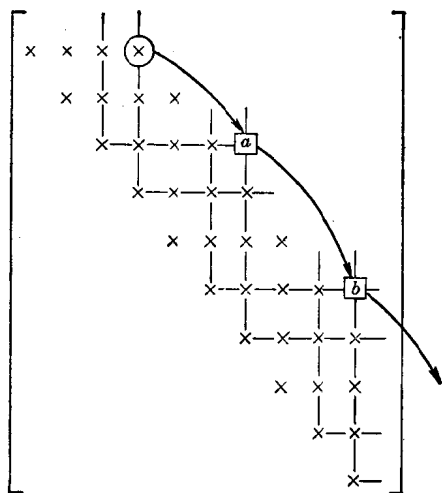


图 2

从上面介绍知道，一个对称矩阵  $A$ ，通过正交变换相似于一个对称三对角矩阵，可以有各种途径，最后得到的相似的对称三对角阵，也可能不相同。但是有下面定理。

**定理 2.3** 对于  $n \times n$  实对称矩阵  $A$ ，如果有正交矩阵  $Q$ ，使得  $Q^T A Q = T$  是一个对称三对角矩阵，并且  $T$  的第一条超对角线上的元素全为正的，那么  $Q$  和  $T$  完全由  $Q$  的第一列  $q_1$  和  $A$  所决定，或者完全由  $Q$  的第  $n$  列  $q_n$  和  $A$  所决定。

**证明** 记

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & 0 \\ \beta_1 & \alpha_2 & \ddots & \\ & \ddots & \ddots & \beta_{n-1} \\ 0 & & \beta_{n-1} & \alpha_n \end{pmatrix},$$

$$Q = (q_1, q_2, \dots, q_n),$$

本定理要证明  $\alpha_i, \beta_i, q_i$  由  $q_1$  和  $A$  唯一确定.

由  $Q^T A Q = T$  得  $AQ = QT$  比较此式两边的各列, 有

$$Aq_1 = \alpha_1 q_1 + \beta_1 q_2, \quad (4)$$

$$Aq_i = \alpha_i q_i + \beta_i q_{i+1} + \beta_{i-1} q_{i-1}, \quad i=2, 3, \dots, n, \quad (5)$$

由  $Q$  是正交阵, 它的各列彼此正交, 因此从 (4) 知

$$\alpha_1 = q_1^T A q_1,$$

$$\beta_1 q_2 = Aq_1 - \alpha_1 q_1,$$

利用  $\beta_1 > 0$  和  $\|q_2\| = 1$ , 故

$$\beta_1 = \|Aq_1 - \alpha_1 q_1\|,$$

从而可知,  $\alpha_1, \beta_1$  完全由  $A$  和  $q_1$  确定. 因为  $\beta_1 > 0$ , 故

$$q_2 = (Aq_1 - \alpha_1 q_1) / \beta_1,$$

也由  $A$  和  $q_1$  完全确定. 一般地如果:  $q_2, q_3, \dots, q_{i-1}, q_i, \alpha_1, \alpha_2, \dots, \alpha_{i-1}, \beta_1, \beta_2, \dots, \beta_{i-1}$  已经由  $A, q_1$  完全确定, 那么利用 (5), 知  $\alpha_i = q_i^T A q_i$ , 若  $i \neq n$  有

$$\beta_i q_{i+1} = Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1},$$

$$\beta_i = \|Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}\|,$$

$$q_{i+1} = (Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}) / \beta_i,$$

完全被  $q_i, q_{i-1}$  和  $A$  所确定, 从而被  $q_1$  和  $A$  所确定. 同理可知所有  $q_1, q_2, \dots, q_{n-1}; \alpha_1, \dots, \alpha_n; \beta_1, \dots, \beta_{n-1}$  完全可由  $A$  和  $q_1$  所确定. 证毕.

定理 2.3 中  $T$  的要求是: 第 1 条超对角线上元素为正. 如果  $T$  的第 1 条超对角线上元素有正有负, 但不为零, 是否

也有类似的结论. 为了回答这个问题先看

**引理 2.1** 若

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & \\ & \ddots & \ddots & \ddots \\ 0 & & \beta_{n-1} & \alpha_n \end{pmatrix}, \quad \beta_i \neq 0, \quad i=1, 2, \dots, n-1,$$

则存在非奇异对角阵  $\Delta = \text{diag}(\delta_1, \delta_2, \dots, \delta_n)$ , 其中  $\delta_1=1$ ,  $\delta_{i+1} = \text{sign}(\beta_i \delta_i)$ ,  $i=1, 2, \dots, n-1$ , 使得

$$\Delta T \Delta^{-1} = \begin{pmatrix} \alpha_1 & |\beta_1| & & 0 \\ |\beta_1| & \alpha_2 & |\beta_2| & \\ & |\beta_2| & \ddots & \ddots \\ 0 & & \ddots & |\beta_{n-1}| \\ & & & \alpha_n \end{pmatrix}. \quad (6)$$

**证明**

$$\begin{aligned} & \Delta T \Delta^{-1} \\ &= \begin{pmatrix} \alpha_1 & (\beta_1 \delta_1 / \delta_2) & & 0 \\ (\beta_1 \delta_2 / \delta_1) & \alpha_2 & (\beta_2 \delta_2 / \delta_3) & \\ & (\beta_2 \delta_3 / \delta_2) & \ddots & \ddots \\ 0 & & \ddots & (\beta_{n-1} \delta_{n-1} / \delta_n) \\ & & & (\beta_{n-1} \delta_n / \delta_{n-1}) & \alpha_n \end{pmatrix}, \end{aligned}$$

因为  $\beta_k \delta_{k+1} / \delta_k = \beta_k \delta_k / \delta_{k+1} = |\beta_k|$ ,

故(6)成立. 证毕.

**定理 2.4** 对于  $n \times n$  实对称阵  $A$ , 如果有正交阵  $P = [p_1, p_2, \dots, p_n]$ , 使得

$$P^T A P = T = \begin{pmatrix} \alpha_1 & \beta_1 & & 0 \\ & \ddots & \ddots & \\ \beta_1 & \alpha_2 & & \beta_{n-1} \\ & & \ddots & \ddots \\ 0 & & \beta_{n-1} & \alpha_n \end{pmatrix}$$

是一个对称三对角矩阵, 并且  $T$  的第 1 条超对角线元素不为 0, 则  $T$  和  $\mathbf{p}_2, \mathbf{p}_3, \dots, \mathbf{p}_n$  完全由  $A$  和  $\mathbf{p}_1$ , 以及  $T$  的第 1 条超对角线上元素的符号所完全确定.

**证明** 如果所有  $\beta_i > 0$ , 则由定理 2.3 知结果成立. 如果  $\beta_i$  中有正有负, 则取引理 2.1 中对角阵  $\Delta$ , 使

$$\Delta T \Delta = T_+ = \begin{pmatrix} \alpha_1 & |\beta_1| & & 0 \\ & \ddots & \ddots & \\ |\beta_1| & \alpha_2 & & |\beta_2| \\ & & \ddots & \ddots \\ 0 & & & |\beta_{n-1}| \\ & & & \ddots & \ddots \\ & & & \beta_{n-1} & \alpha_n \end{pmatrix},$$

于是

$$\Delta P^T A P \Delta^{-1} = \Delta T \Delta^{-1} = T_+,$$

记  $Q = P \Delta^{-1} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n) = \left( \mathbf{p}_1, \frac{1}{\delta_2} \mathbf{p}_2, \frac{1}{\delta_3} \mathbf{p}_3, \dots, \frac{1}{\delta_n} \mathbf{p}_n \right)$ ,

$Q^T Q = \Delta^{-1} P^T P \Delta^{-1} = \Delta^{-1} \Delta^{-1} = I$ , 故  $Q$  是一个正交阵. 由定理 2.3 知  $Q$  的各列和  $T_+$  完全被  $A$  和  $\mathbf{q}_1$  所确定, 而  $\mathbf{q}_1 = \mathbf{p}_1$ , 因此  $Q$  的各列和  $T_+$  完全被  $A$  和  $\mathbf{p}_1$  所确定. 知道了  $\beta_i$  的符号, 就可以从  $|\beta_i|$  得到  $\beta_i$ , 从  $\delta_i \mathbf{q}_i$  得到  $\mathbf{p}_i$ ; 故  $A$  和  $\mathbf{p}_1$  完全确定  $P$  和  $T$ . 证毕.

本节所介绍的通过镜像相似变换, 将  $A$  化成对称三对角矩阵, 是将  $A$  左乘和右乘一系列初等阵  $Q_i = I - 2\mathbf{w}_i \mathbf{w}_i^T$  后, 相似成一个对称三对角阵  $T_1$ , 即

$$T_1 = Q_s Q_{s-1} \cdots Q_2 Q_1 A Q_1 Q_2 \cdots Q_s,$$

$Q_1 Q_2 \cdots Q_s$  是一个正交阵. 注意到  $w_i$  的第 1 个分量总是零, 因此  $Q_i$  的第 1 列为  $e_1 = (1, 0, \dots, 0)^T$ , 第 1 行为  $e_1^T = (1, 0, 0, \dots, 0)$ , 即

$$Q_i = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & & & & \\ \vdots & & H_i & & \\ 0 & & & & \end{pmatrix},$$

因此  $Q = Q_1 Q_2 \cdots Q_s$  的第 1 列为  $e_1$ ; 如果镜像变换结果的第一条超对角线上元素  $\tau$  都不为 0, 那么  $T_1$  和  $Q$  完全由  $A$  和  $e_1$  以及那些  $\tau$  的符号所决定. 因此这个过程与定理 2.3 中的下述过程产生的结果  $T$  和  $Q$  是一样的:

$$\begin{aligned} \text{取} \quad q_1 &= e_1, \quad \beta_0 = 0, \\ \alpha_i &= (Aq_i, q_i), \\ \beta_i q_{i+1} &= Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}, \\ \beta_i &= \pm \|Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}\| \\ &\quad (\text{符号由对应 } \tau \text{ 的符号确定}), \\ q_{i+1} &= (Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}) / \beta_i, \\ &\quad i = 1, 2, \dots, n-1, \\ \alpha_n &= (Aq_n, q_n). \end{aligned}$$

由此我们也可以知道, 镜像变换过程中, 出现  $\tau=0$  的充分必要条件是  $n$  个向量

$$e_1, Ae_1, \dots, A^{n-1}e_1$$

线性相关. 这是因为  $\beta_i=0$ , 充要条件是  $Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1} = 0$ , 而后者意味着

$$e_i, Ae_1, \dots, A^{i-1}e_1, A^i e_1$$

线性相关.

再来考察平面旋转变换将  $A$  化成对称三对角矩阵的过

程. 记第  $k$  次旋转时, 使用的  $s, c$  为  $s_k, c_k$ ,  $t = \frac{(n-1)(n-2)}{2}$ .

设  $A$  是满秩矩阵, 按足标次序 (3), 通过平面旋转, 将  $A$  相似于对称三对角阵  $T_2$ , 于是

$$T_2 = T_{n-1,n}(s_t, c_t) \cdots T_{n-2,n-1}(s_2, c_2) T_{n-1,n}(s_1, c_1) A \\ T_{n-1,n}^T(s_1, c_1) T_{n-2,n-1}^T(s_2, c_2) \cdots T_{n-1,n}^T(s_t, c_t),$$

记  $P = T_{n-1,n}^T(s_1, c_1) T_{n-2,n-1}^T(s_2, c_2) \cdots T_{n-1,n}^T(s_t, c_t)$ , 在这些旋转矩阵  $T_{j-1,j}^T$  中, 都有  $j \geq 3$ , 因此  $T_{j-1,j}^T(s_k, c_k)$  的第 1 行第 1 列都与单位阵  $I$  的第 1 行第 1 列相同. 这样也可知道  $P$  的第 1 列为  $e_1$ , 于是可知若  $T_1$  和  $T_2$  的第 1 条超对角线上的元素全为正, 则  $T_2 = T_1$ ,  $P = Q$ . 如果  $T_2$  的第 1 条超对角线上的元素有正有负, 但不为零, 则  $T_1 = \Delta T_2 \Delta^{-1}$ ,  $Q = P \Delta^{-1}$ , 这里  $\Delta = \text{diag}(1, \delta_2, \cdots, \delta_n)$  是象引理 2.1 中所给的  $\Delta$  那样的矩阵.

上面指出镜像变换和旋转变换, 都是  $q_1 = e_1$  的正交相似变换

$$Q^T A Q = T,$$

将  $A$  变成对称三对角阵  $T$ . 当然也可以取任意单位向量  $a$  为  $q_1$ , 此时进行的过程为

$$\begin{cases} q_1 = a, \\ \alpha_i = (Aq_i, q_i), \beta_i q_{i+1} = Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}, \\ \beta_i = \|Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}\|, \\ q_{i+1} = (Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}) / \beta_i, \\ i = 1, 2, \cdots, n-1, \end{cases} \quad (7)$$

(7) 称为 Lanczos 过程.

如果在过程中某个  $\beta_i = 0$ , 则可取任何与  $q_1, q_2, \cdots, q_i$  都正交的单位向量作为  $q_{i+1}$ , 过程可继续进行.



### §3 对称三对角矩阵的极值性质\*

记对称三对角矩阵

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \ddots & \ddots & \ddots & \\ 0 & & \ddots & \beta_{n-1} & \alpha_n \end{pmatrix}$$

的截段

$$\begin{pmatrix} \alpha_l & \beta_l & & & 0 \\ \beta_l & \alpha_{l+1} & \beta_{l+1} & & \\ & \ddots & \ddots & \ddots & \\ 0 & & \ddots & \beta_{m-l} & \alpha_m \end{pmatrix}$$

为  $T_{l,m}$ ,  $1 \leq l < m \leq n$ ; 它的特征多项式为  $\chi_{l,m}(\lambda)$ , 因此  $T_{1,n} = T$ . 为了方便起见记  $T_j = T_{1,j}$ , 特征多项式为  $\chi_j(\lambda)$ .

**引理 2.2** 若  $l < j$  则  $T_j^l e_1$  的第  $l+1$  个分量为  $\beta_1 \beta_2 \cdots \beta_l$ , 而第  $l+2$  个分量至第  $j$  个分量都为 0, 这里  $e_1$  是  $j$  维单位向量  $(1, 0, \dots, 0)^T$ .

**证明**

$$T_j e_1 = \begin{pmatrix} \alpha_1 \\ \beta_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

\* 本节引理 2.2, 定理 2.5 的结果取自[10].

第2个分量为 $\beta_1$ , 第3个分量至第 $j$ 个分量都为0, 因此命题对 $l=1$ 成立. 今假设命题对 $l=k < j-1$ 成立.

$T_j^k \mathbf{e}_1$ 的第 $k+1$ 个分量为 $\beta_1 \beta_2 \cdots \beta_k$ , 而第 $k+2$ 个分量至第 $j$ 个分量都为0. 考察 $T_j T_j^k \mathbf{e}_1$ 的第 $k+2$ 个元素. 因为 $T_j$ 的第 $k+2$ 行为三个元素 $\beta_{k+1}, \alpha_{k+2}, \beta_{k+2}$ 分别在第 $k+1$ 列, 第 $k+2$ 列, 第 $k+3$ 列, 因此 $T_j T_j^k \mathbf{e}_1$ 的第 $k+2$ 个元素为 $\beta_1 \beta_2 \cdots \beta_k \beta_{k+1}$ . 而 $T_j$ 的第 $k+3$ 行的三个元素 $\beta_{k+2}, \alpha_{k+3}, \beta_{k+3}$ 是在第 $k+2$ 列,  $k+3$ 列,  $k+4$ 列, 因此 $T_j T_j^k \mathbf{e}_1$ 的第 $k+3$ 个元素为0, 同样可以证明次序数大于 $k+3$ 的元素也为0. 证毕.

**定理 2.5** 设 $T$ 的第1条超对角线上的元素全为正, 则对 $j=1, 2, \dots, n-1$ 成立

$$a. \beta_1 \beta_2 \cdots \beta_j = \|\chi_j(T) \mathbf{e}_1\| = \min_{\psi \in mP_j} \|\psi(T) \mathbf{e}_1\|,$$

这里 $mP_j$ 表示全体首项系数为1的 $j$ 次多项式的集合.

$$b. \mathbf{e}_{j+1} = \chi_j(T) \mathbf{e}_1 / \beta_1 \beta_2 \cdots \beta_j.$$

**证明** 由引理2.2知对任意 $\psi(\lambda) \in mP_j$ ,  $\psi(T)$ 的第 $j+1$ 个分量为 $\beta_1 \beta_2 \cdots \beta_j$ , 因此

$$\beta_1 \beta_2 \cdots \beta_j \leq \|\psi(T) \mathbf{e}_1\|,$$

特别取 $mP_j$ 中的 $\chi_j(\lambda)$ , 证明

$$\chi_j(T) \mathbf{e}_1 = \beta_1 \beta_2 \cdots \beta_j \mathbf{e}_{j+1}.$$

记 $E_j = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_j)$ 是 $n \times j$ 矩阵. 首先证明下列命题

$$E_j^* T^l \mathbf{e}_1 = T_j^l \mathbf{e}_1$$

对 $l \leq j$ 成立. 实际上此命题意思即为,  $T^l \mathbf{e}_1$ 的为首 $j$ 个分量与 $T_j^l \mathbf{e}_1$ 的分量完全相同. 对于 $l=1, j=1$ 这命题显然成立, 对于 $l=1, j \geq 2$

$$T\mathbf{e}_1 = \begin{pmatrix} \alpha_1 \\ \beta_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{1 \times n}, \quad E_j^* T\mathbf{e}_1 = \begin{pmatrix} \alpha_1 \\ \beta_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{1 \times j} = T_j \mathbf{e}_1,$$

因此命题对  $l=1$  成立. 假定它对  $k \leq j-1$  也成立. 考虑  $T^{k+1}\mathbf{e}_1 = T(T^k\mathbf{e}_1)$ , 因为  $T^k\mathbf{e}_1$  的第  $k+2$  个分量开始后面的分量全为 0, 因此第  $j+1$  个分量以及后面的分量都为 0, 这样  $T(T^k\mathbf{e}_1)$  的为首  $j$  个分量与  $T_j(T_j^k\mathbf{e}_1)$  的相同, 命题得证.

现在来讨论  $\chi_j(T)\mathbf{e}_1$ , 由引理 2.2 知道它的第  $j+1$  个分量为  $\beta_1\beta_2\cdots\beta_j$ , 另外  $\chi_j(T)\mathbf{e}_1$  的为首  $j$  个分量与  $\chi_j(T_j)\mathbf{e}_1$  的  $j$  个分量相同, 而后者根据 Cayley-Hamilton 定理(见[9]中 p. 363)知  $\chi_j(T_j)=0$ , 即  $\chi_j(T_j)\mathbf{e}_1=0$ , 因此  $\chi_j(T)\mathbf{e}_1$  除了第  $j+1$  个分量外全为 0, 即

$$\chi_j(T)\mathbf{e}_1 = \beta_1\beta_2\cdots\beta_j\mathbf{e}_{j+1}. \text{ 证毕.}$$

**推论** 设  $Q=[\mathbf{q}_1\mathbf{q}_2\cdots\mathbf{q}_n]$  是一个正交矩阵, 又设  $A$  是实对称矩阵. 若  $Q^*AQ=T$  是一个对称三对角矩阵,  $\beta_i>0$ ,  $i=1, 2, \dots, n-1$ , 则

$$\beta_1\beta_2\cdots\beta_j = \|\chi_j(A)\mathbf{q}_1\| = \min_{\psi \in mP_j} \|\psi(A)\mathbf{q}_1\| \quad (8)$$

$$\text{且} \quad \mathbf{q}_{j+1} = \chi_j(A)\mathbf{q}_1 / \beta_1\beta_2\cdots\beta_j. \quad (9)$$

**证明** 由  $Q^*AQ=T$ , 因此  $Q^*\chi_j(A)Q = \chi_j(T)$ , 故

$$\chi_j(T)\mathbf{e}_1 = Q^*\chi_j(A)Q\mathbf{e}_1 = Q^*\chi_j(A)\mathbf{q}_1, \quad (10)$$

由定理 2.5

$$\begin{aligned} \beta_1\beta_2\cdots\beta_j &= \|\chi_j(T)\mathbf{e}_1\| = \|\chi_j(A)\mathbf{q}_1\| \\ &\leq \min_{\psi \in mP_j} \|\psi(T)\mathbf{e}_1\| = \min_{\psi \in mP_j} \|\psi(A)\mathbf{q}_1\|, \end{aligned}$$

(8)式得到证明, 另外从(10)式知

$$\beta_1 \beta_2 \cdots \beta_j e_{j+1} = Q^* \chi_i(A) q_1,$$

两边左乘  $Q$ , 即得(9)式, 从而本推论证毕.

在对称方程组求解问题, 和对称矩阵特征值问题中, 我们常碰到如下的线性方程组

$$T_j f_j = \alpha e_1,$$

这里  $\alpha$  是常数,  $f_j = (\delta_1^{(j)}, \delta_2^{(j)}, \dots, \delta_j^{(j)})^T$  是未知向量, 要求将  $\delta_i^{(j)}$  用  $T_j$  的元素和  $\alpha$  表示出来. 有如下定理(参见[31])

**定理 2.6** 若  $\det(T_j) = d_j \neq 0$ , 则

$$\delta_k^{(j)} = (-1)^{k-1} \alpha \frac{\beta_1 \beta_2 \cdots \beta_{k-1} \det(T_{k+1,j})}{\det(T_j)}, \quad (11)$$

$$k=1, 2, \dots, j, \quad \text{这里} \quad \det(T_{j+1,j}) = 1.$$

**证明** 利用 Cramer 法则

$$\delta_k^{(j)} = s_k / d_j, \quad (12)$$

这里

$$s_k = \det \begin{pmatrix} \alpha_1 & \beta_1 & & & \alpha & & & & \\ & \alpha_2 & \ddots & & \vdots & & & & \\ & \beta_1 & \alpha_2 & \ddots & 0 & & & & \\ & & \ddots & \ddots & 0 & & & & \\ & & & \beta_{k-2} & \alpha_{k-1} & 0 & & & \\ & & & & \beta_{k-1} & 0 & \beta_k & & \\ & & & & & 0 & \alpha_{k+1} & \beta_{k+2} & \\ & & & & & & \beta_{k+1} & \alpha_{k+2} & \ddots \\ & & & & & & & \ddots & \ddots \\ & & & & & & & & 0 & \beta_{j-1} & \alpha_j \end{pmatrix},$$

将  $s_k$  按第  $k$  列展开,



1968 年 Thompson 和 McEntegert[12] 给出如下的特征向量与伴随阵之间的关系:

**定理 2.7** 设  $n \times n$  实对称矩阵  $A$  的  $n$  个标准正交特征向量为  $z_1, z_2, \dots, z_n$ , 它们对应的特征值依次为  $\lambda_1, \lambda_2, \dots, \lambda_n$ ,  $\chi(\lambda)$  为  $A$  的特征多项式  $\det(\lambda I - A)$ , 则

$$\text{adj}(\lambda_i I - A) = \chi'(\lambda_i) z_i z_i^*.$$

**证明** 取常数  $\mu \neq \lambda_j, j=1, 2, \dots, n$ , 于是  $\mu I - A$  有逆

$$\text{adj}(\mu I - A) = \det(\mu I - A) (\mu I - A)^{-1},$$

记  $Z = [z_1 z_2 \dots z_n]$  是一个正交阵,  $\Delta = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ , 有

$$A = Z \Delta Z^*,$$

$$\text{故} \quad \text{adj}(\mu I - A) = \chi(\mu) Z (\mu I - \Delta)^{-1} Z^*,$$

或

$$\text{adj}(\mu I - A) = Z \Delta(\mu) Z^*, \quad (13)$$

$$\text{其中} \quad \Delta(\mu) = \text{diag}\left(\frac{\chi(\mu)}{\mu - \lambda_1}, \frac{\chi(\mu)}{\mu - \lambda_2}, \dots, \frac{\chi(\mu)}{\mu - \lambda_n}\right),$$

等式(13)的两边都是  $\mu$  的解析函数, 因此对任何复数都成立. 特别取  $\mu = \lambda_i$ ,

$$\Delta(\mu) = \Delta(\lambda_i) = \text{diag}(\overset{\text{第 } i \text{ 列}}{0}, 0, \dots, \chi'(\lambda_i), 0, \dots, 0),$$

于是  $\text{adj}(\lambda_i I - A) = Z \Delta(\lambda_i) Z^* = \chi'(\lambda_i) z_i z_i^*$ . 证毕.

**定理 2.8** 将特征向量的分量与伴随阵元素联系起来. 因此对于形式简单的矩阵, 就有可能将特征向量的分量与矩阵的元素联系起来. 1971 年 C. Paige 就给出了如下的定理 [13].

**定理 2.8** 设对称三对角矩阵  $T$  的特征值为  $\theta_1, \theta_2, \dots, \theta_n$ , 对应的标准化的正交特征向量为  $s_1, s_2, \dots, s_n$ , 则对  $\mu \leq \nu$ ,  $j=1, 2, \dots, n$  有

$$\chi'_{1,n}(\theta_j) s_{\mu j} s_{\nu j} = \chi_{1,\mu-1}(\theta_j) \beta_{\mu} \dots \beta_{\nu-1} \chi_{\nu+1,n}(\theta_j), \quad (14)$$

这里  $s_{\mu j}$  是  $\mathbf{s}_j$  的第  $\mu$  个分量. 特别当  $T$  不可约时 (即  $\beta_i \neq 0$ ,  $i=1, 2, \dots, n-1$ ) 有

$$s_{\mu j}^2 = \chi_{1, \mu-1}(\theta_j) \chi_{\mu+1, n}(\theta_j) / \chi'_{1, n}(\theta_j). \quad (15)$$

证明 先证(14)式. 利用定理 2.8 有

$$\text{adj}(\theta_j I - T) = \chi'_{1, n}(\theta_j) \mathbf{s}_j \mathbf{s}_j^*, \quad (16)$$

等式右边的第  $\mu$  行第  $\nu$  列即为  $\chi'_{1, n}(\theta_j) s_{\mu j} s_{\nu j}$ ; 等式左边的第  $\mu$  行第  $\nu$  列是  $\theta_j I - T$  的第  $\nu$  行第  $\mu$  列的代数余子式, 因为  $\theta_j I - T$  是对称的, 因此也是  $\theta_j I - T$  的第  $\mu$  行第  $\nu$  列的代数余子式. 这个代数余子式记为  $\Delta_{\mu, \nu}$ .

为了求  $\Delta_{\mu, \nu}$  我们将  $\theta_j I - T$  划去第  $\mu$  行和第  $\nu$  列后的矩阵表示如下:

$$\begin{pmatrix} \theta_j I - T_{1, \mu-1} & 0 & 0 & 0 \\ & -\beta_{\mu-1} & & \\ & 0 & \begin{pmatrix} -\beta_{\mu} & & 0 \\ & -\beta_{\mu+1} & \dots & 0 \\ & 0 & \dots & -\beta_{\nu-1} \end{pmatrix} & 0 \\ & 0 & 0 & \theta_j I - T_{\nu+1, n} \end{pmatrix},$$

因此

$$\begin{aligned} \Delta_{\mu \nu} &= (-1)^{\mu+\nu} \det(\theta_j I - T_{1, \mu-1}) (-1)^{\nu-\mu} \\ &\quad \times \beta_{\mu} \beta_{\mu+1} \cdots \beta_{\nu-1} \det(\theta_j I - T_{\nu+1, n}) \\ &= \chi_{1, \mu-1}(\theta_j) \beta_{\mu} \beta_{\mu+1} \cdots \beta_{\nu-1} \chi_{\nu+1, n}(\theta_j), \end{aligned}$$

(14)式获证.

当  $\nu = \mu$  时, (16) 式的右边第  $\mu$  行第  $\mu$  列为  $\chi'_{1, n}(\theta_j) s_{\mu j}^2$ , 而左边的第  $\mu$  行第  $\mu$  列为  $\chi_{1, \mu-1}(\theta_j) \chi_{\mu+1, n}(\theta_j)$ . 又因为  $T$  不可约时它没有重特征值, 因此  $\chi'_{1, n}(\theta_j) \neq 0$ . 故(15)式成立. 证毕.

## 推论

$$s_{1j}s_{nj}\chi'(\theta_j) = \beta_1\beta_2\cdots\beta_{n-1}, \quad (17)$$

$$s_{1j}^2\chi'(\theta_j) = \chi_{2,n}(\theta_j), \quad (18)$$

$$s_{nj}^2\chi'(\theta_j) = \chi_{1,n-1}(\theta_j), \quad (19)$$

这里记  $\chi(\theta) = \chi_{1,n}(\theta)$ .

**证明** 在定理中取  $\mu=1$ ,  $\nu=n$  即得(17), 取  $\mu=\nu=1$  即得(18), 取  $\mu=\nu=n$  得到(19). 证毕.

从(17)式可知, 如果  $T$  不可约, 那么  $\beta_1\beta_2\cdots\beta_{n-1} \neq 0$ ,  $\chi'(\theta_j) \neq 0$ , 因此  $s_{1j}s_{nj} \neq 0$ , 即任意一个特征向量的第 1 个分量和最后一个分量不为 0. 再由(18)、(19)两式知  $\chi_{2,n}(\theta_j) \neq 0$ ,  $\chi_{1,n-1}(\theta_j) \neq 0$ . 这也可从 Jacobi 矩阵根的分隔性质而得到.

现在来考虑矩阵的特征值反问题, 也即从已知  $n$  个特征值来确定一个  $n \times n$  矩阵. 当然  $n \times n$  矩阵有  $n^2$  个元素, 一般情况不能被  $n$  个特征值所确定. 因此还要提供另外的一些信息, 才能完全确定一个矩阵.

当  $n \times n$  矩阵是一个对称三对角矩阵时, 要确定的矩阵元素只有  $2n-1$  个, 因此需要附加的信息也较简单.

1967 年 Harry Hochstadt[14] 提出二种类型问题.

**问题 1** 给定二个序列  $\lambda_1, \lambda_2, \dots, \lambda_n$  和  $\mu_1, \mu_2, \dots, \mu_{n-1}$  并且  $\lambda_i < \mu_i < \lambda_{i+1}$ ,  $i=1, 2, \dots, n-1$ , 要求构造一个  $n \times n$  对称三对角矩阵  $T_{1,n}$ ,  $\beta_i > 0$ , 使得  $\lambda_1, \lambda_2, \dots, \lambda_n$  是  $T_{1,n}$  的特征值, 而  $\mu_1, \mu_2, \dots, \mu_{n-1}$  是  $T_{2,n}$  的特征值.

**问题 2** 给定  $n$  个由小到大的数  $\lambda_1, \lambda_2, \dots, \lambda_n$ , 要求构造一个  $n \times n$  对称三对角阵  $T$ , 使得  $\lambda_1, \lambda_2, \dots, \lambda_n$  是  $T$  的特征值, 但是  $T$  的元素满足:

$$\alpha_i = \alpha_{n+1-i},$$

$$\beta_i = \beta_{n-i} > 0,$$



$$i=1, 2, \dots, \left[\frac{n}{2}\right],$$

也即  $T$  是关于第二条主对角线(从右上到左下)也是对称的。

Hochstadt 证明问题 1 和问题 2 都至多有一个解。1976 年 Ole, H. Hald 证明问题 1 和问题 2 都至少有一个解, 并且给出求解的方法 [15], 但是算法不稳定, 1978 年 O. de Boor 和 G. H. Golub, 给出一种数值稳定的算法来求上述问题 1 和问题 2 的解 [16], 我们下面介绍的 B. N. Parlett 教授在 [10] 中给出的, 利用 C. Paige 公式来构造  $T$  的方法。在介绍这项方法之前, 先来简单介绍一下对称三对角矩阵的特征值反问题实际背景。

考虑一个两端固定的弦的振动问题, 它的数学模型可以化成下列边值问题:

$$\begin{cases} u''(x) - \sigma(x)u(x) = \lambda u(x); \\ u(0) = u(1) = 0, \end{cases} \quad (20)$$

$\sigma(x)$  是跟弦的密度有关的函数, 如果弦是均匀的, 那么  $\sigma(x) = \sigma_0 = \text{const}$ . 当  $\lambda$  取固有值时, (20) 有非零解  $u(x)$ , 这是固有值对应弦的固有频率。

已知  $\sigma(x)$ , 求固有值  $\lambda$  的问题, 称为固有值问题或特征值问题。反过来, 如果通过实验手段测得弦的固有频率, 计算得到固有值, 从已知固有值来求  $\sigma(x)$  的问题, 就称为特征值反问题。

用差分方法解边值问题 (20), 就把微分方程的特征值问题, 化成矩阵的特征值问题, 因此也把微分方程的特征值反问题, 化成矩阵的特征值反问题。例如  $u''(x)$  用中心差分  $u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))/h^2$  代替, 这里  $h=1/n$ ,  $x_i = ih$ , 那么对应的矩阵是一个对称三角阵。因此弦振动问题的特征值

反问题,就化成一个对称三对角矩阵的特征值反问题.

边值问题(20)是弦振动问题的数学模型,也可看作一个粒子的一维薛定谔方程,这时  $u(x)$  表示粒子的波函数,  $\sigma(x)$  代表势函数,固有值  $\lambda$  代表在这个场中运动的粒子所具有的能级. 如果已知各个允许的能级,求势函数  $\sigma(x)$  的问题也是对称三对角矩阵的特征值反问题.

上面介绍的简单的物理问题对应对称三对角矩阵的特征值反问题,可以设想,更复杂一些的某些物理问题,会对应更复杂一些矩阵的特征值反问题. 因此从七十年代以来矩阵特征值反问题的研究,被专家们所重视.

现在回到问题 1、2 的解上来. 设  $T$  是对称三对角矩阵,  $\theta_1, \theta_2, \dots, \theta_n$  是它的特征值,  $s_1, s_2, \dots, s_n$  是对应的特征向量,  $S = [s_1, s_2, \dots, s_n]$  是正交阵,  $\Lambda = \text{diag}(\theta_1, \theta_2, \dots, \theta_n)$ , 于是有

$$S^*TS = \Lambda,$$

$$T = SAS^*,$$

将  $\Lambda$  看作一个对称矩阵  $A$ , 利用 § 2 的唯一归化定理知  $T$  和  $S^*$  的其他各列,完全被  $\Lambda$  和  $S^*$  的第 1 列  $(s_{11}, s_{12}, \dots, s_{1n})^T$  所确定,这里  $s_{1i}$  是  $T$  的第  $i$  个特征向量的第一个分量. 因此如果由  $\Lambda$  和反问题所给的附加信息能够确定  $(s_{11}, s_{12}, \dots, s_{1n})$ , 那么就能完全确定  $T$ .

对于问题 1, 已知  $\mu_1, \mu_2, \dots, \mu_{n-1}$  是  $T_{2,n}$  的特征值, 因此  $\chi_{2n}(\lambda) = \det(\lambda I - T_{2,n}) = \prod_{i=1}^{n-1} (\lambda - \mu_i)$ , 由 (18)

$$s_{1j}^2 \chi'(\theta_j) = \chi_{2n}(\theta_j),$$

即得 
$$s_{1j}^2 = \prod_{i=1}^{n-1} (\theta_j - \mu_i) / \prod_{i=1, i \neq j}^n (\theta_j - \theta_i),$$

因为特征向量可以差一个符号，可取

$$s_{ij} = \sqrt{\prod_{i=1}^{n-1} (\theta_j - \mu_i)} / \prod_{\substack{i=1 \\ i \neq j}}^n (\theta_j - \theta_i),$$

再由 Lanczos 过程(7)，就可以得到  $T$ 。

对于问题 2，因为  $T$  是关于第 2 条主对角线对称的，

$$\alpha_i = \alpha_{n+1-i}, \quad \beta_i = \beta_{n-i} > 0,$$

于是若记矩阵  $\tilde{I} = (e_n, e_{n-1}, \dots, e_2, e_1)$ ，有

$$T = \tilde{I} T \tilde{I},$$

由

$$T s_i = \lambda_i s_i,$$

$$\tilde{I} T \tilde{I} s_i = \lambda_i s_i,$$

$$T \tilde{I} s_i = \lambda_i \tilde{I} s_i,$$

因为  $\beta_i > 0$ ， $T$  不可约， $T$  的特征值都不相同，因此  $\tilde{I} s_i$  与  $s_i$  至多差一个符号，即

$$\tilde{I} s_i = \varepsilon_i s_i, \quad \varepsilon_i = \pm 1,$$

由此可知

$$s_m = \varepsilon_i s_{1i},$$

利用 (17)

$$s_{1i} s_{ni} \chi'(\theta_i) = \beta_1 \beta_2 \cdots \beta_{n-1},$$

或

$$s_{1i}^2 = \frac{\beta_1 \beta_2 \cdots \beta_{n-1}}{\varepsilon_i \chi'(\theta_i)},$$

但  $\chi'(\theta_i) = (\theta_i - \theta_1) \cdots (\theta_i - \theta_{i-1}) (\theta_i - \theta_{i+1}) \cdots (\theta_i - \theta_n)$ ,

$\text{sign } \chi'(\theta_i) = (-1)^{n-i}$ ， $\beta_1 \cdots \beta_{n-1} > 0$ ，故  $\varepsilon_i = (-1)^{n-i}$ ，

从而得到

$$s_{1i} = \sqrt{\beta_1 \beta_2 \cdots \beta_{n-1}} / \prod_{l=1}^{i-1} (\theta_i - \theta_l) \prod_{l=i+1}^n (\theta_l - \theta_i),$$

由此也可求得  $T$ 。

现在我们把上面的理论，应用到如下的一个有实际意义的问题中去。考虑如图所示的  $n$  个质点， $n$  段弹簧的振动模

型,  $m_i$  表示第  $i$  个质点的质量,  $k_i$  表示第  $i$  段弹簧的弹性系数. 若已知这个弹性系统的全部固有频率, 又知道这个弹性系统的最后一段弹簧和最后一个质点去掉后的弹性系统的全部固有频率; 同时知道  $m_1$ , 要求出  $k_1, k_2, \dots, k_n$  和  $m_2, m_3, \dots, m_n$  的值.

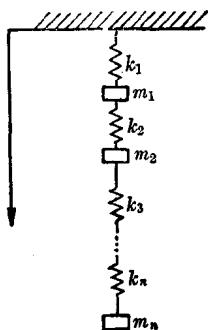


图 3

要解这个问题, 先把振动方程列出来. 设  $l_i$  表示第  $i$  段弹簧在静止状态时的长度,  $u_i$  表示第  $i$  段弹簧的伸长,  $x_i$  表示第  $i$  个质点的坐标. 于是

$$x_i = \sum_{j=1}^i u_j + \sum_{j=1}^i l_j,$$

按照牛顿力学方程, 对第  $i$  个质点成立

$$m_i \ddot{x}_i = k_{i+1} u_{i+1} - k_i u_i,$$

即

$$m_i \sum_{j=1}^i \ddot{u}_j = k_{i+1} u_{i+1} - k_i u_i,$$

或

$$\sum_{j=1}^i \ddot{u}_j = \frac{k_{i+1}}{m_i} u_{i+1} - \frac{k_i}{m_i} u_i, \quad i=1, 2, \dots, n, \quad (21)$$

$$k_{n+1} = 0,$$

从(21)可得

$$\begin{aligned} \ddot{u}_i &= \frac{k_{i-1}}{m_{i-1}} u_{i-1} - \left( \frac{k_i}{m_{i-1}} + \frac{k_i}{m_i} \right) u_i \\ &\quad + \frac{k_{i+1}}{m_i} u_{i+1}, \quad i=1, 2, \dots, n, \end{aligned} \quad (22)$$

其中规定  $\frac{k_0}{m_0} = \frac{k_1}{m_0} = 0$ .

若记  $a_i = -\left( \frac{k_i}{m_{i-1}} + \frac{k_i}{m_i} \right)$ ,  $b_i = \frac{k_i}{m_i}$ ,  $c_i = \frac{k_{i+1}}{m_i}$ ,

$$A_{1,n} = \begin{pmatrix} a_1 & c_1 & & 0 \\ b_1 & a_2 & c_2 & \\ & b_2 & \ddots & c_{n-1} \\ 0 & & \ddots & b_{n-1} & a_n \end{pmatrix},$$

$\mathbf{u} = (u_1, u_2, \dots, u_n)^T$ , 则方程(22)可表示成

$$\ddot{\mathbf{u}} = A_{1,n} \mathbf{u},$$

若  $A_{1,n}$  的特征值为  $-\lambda_l^2$ ,  $l=1, 2, \dots, n$ , 则  $\lambda_l$ ,  $l=1, 2, \dots, n$ ,

即为这个系统的全部固有频率。同时若

$$A_{1,n-1} = \begin{pmatrix} a_1 & c_1 & & 0 \\ b_1 & a_2 & c_2 & \\ & b_2 & \ddots & c_{n-2} \\ 0 & & \ddots & b_{n-2} & a_{n-1} \end{pmatrix}$$

的特征值为  $-\mu_1^2, -\mu_2^2, \dots, -\mu_{n-1}^2$ , 则  $\mu_1, \mu_2, \dots, \mu_{n-1}$  即为去掉弹簧  $k_n$  和质量  $m_n$  后的系统的全部固有频率。

为了应用本节的理论, 从  $\lambda_l$  和  $\mu_l$  求  $k_i, m_i$ , 我们先要将  $A_{1,n}$  对称化。若

$$D_n = \text{diag}(1, \delta_1, \delta_2, \dots, \delta_{n-1}),$$

有

$$D_n^{-1} A_{1,n-1} D_n$$

$$= \begin{pmatrix} a_1 & c_1 \delta_1 & & 0 \\ b_1 / \delta_1 & a_2 & c_2 \delta_2 / \delta_1 & \\ & b_2 \delta_1 / \delta_2 & \ddots & c_{n-1} \delta_{n-1} / \delta_{n-2} \\ 0 & & \ddots & b_{n-1} \delta_{n-2} / \delta_{n-1} & a_n \end{pmatrix},$$

如果取  $\delta_1, \delta_2, \dots, \delta_{n-1}$ , 使得

$$b_i \delta_{i-1} / \delta_i = c_i \delta_i / \delta_{i-1}, \quad i=1, 2, \dots, n-1,$$

就能使

$$B_{1,n} = D_n^{-1} A_{1,n} D_n$$

是一个对称三对角矩阵, 此时

$$\delta_i = (b_1 b_2 \dots b_i / c_1 c_2 \dots c_i)^{\frac{1}{2}},$$

$$B_{1,n} = \begin{pmatrix} a_1 & \beta_1 & & & 0 \\ \beta_1 & a_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \\ 0 & & \ddots & \beta_{n-1} & a_n \\ & & & \beta_{n-1} & a_n \end{pmatrix},$$

其中  $\beta_i = \sqrt{b_i c_i} = \sqrt{k_i k_{i+1}} / m_i$ .

同时可知  $B_{1,n-1}$  相似于  $A_{1,n-1}$ . 这样我们的问题化成已知  $B_{1,n}$  和  $B_{1,n-1}$  的特征值, 以及  $m_1$  求  $m_2, \dots, m_n$  和  $k_1, \dots, k_n$ . 如果知道了  $B_{1,n}$  的全部元素  $a_1, a_2, \dots, a_n$  和  $\beta_1, \beta_2, \dots, \beta_{n-1}$ , 那么可以从如下递推关系式

$$\begin{aligned} k_1 &= -a_1 m_1, \\ k_i &= \beta_{i-1}^2 m_{i-1}^2 / k_{i-1}, \\ m_i &= -k_i / (a_i + k_i / m_{i-1}), \quad i=2, 3, \dots, n, \end{aligned} \quad (23)$$

求出  $k_i$  和  $m_i$ .

由此问题化成已知  $B_{1,n}$  和  $B_{1,n-1}$  的特征值求  $B_{1,n}$  的全部元素. 这就是 Hochstadt 的问题 1. 不过我们现在知道的是  $B_{1,n-1}$  的特征值而不是  $B_{2,n}$  的特征值.

如果  $s = (s_1, s_2, \dots, s_n)$  是  $B_{1,n}$  特征向量构成的正交阵,

$$\Lambda = \text{diag}(-\lambda_1^2, -\lambda_2^2, \dots, -\lambda_n^2),$$

由

$$B_{1,n} = S \Lambda S^*,$$

可知  $B_{1,n}$  完全被  $S^*$  的最后一列和  $\Lambda$  所完全确定. 但  $S^*$  的

最后一列为

$$\mathbf{g}_n = (s_{n1}, s_{n2}, \dots, s_{nn})^T,$$

而利用公式(19),

$$s_{nj}^2 \chi'(-\lambda_j^2) = \chi_{1,n-1}(-\lambda_j^2),$$

可以求出  $(s_{n1}, s_{n2}, \dots, s_{nn})^T$ , 这里  $\chi(\lambda)$  和  $\chi_{1,n-1}(\lambda)$  分别是  $B_{1,n}$  和  $B_{1,n-1}$  的特征多项式, 即

$$s_{nj}^2 = \chi_{1,n-1}(-\lambda_j^2) / \chi'(-\lambda_j^2) = \prod_{i=1}^{n-1} (\mu_i^2 - \lambda_j^2) / \prod_{\substack{i=1 \\ i \neq j}}^n (\lambda_i^2 - \lambda_j^2), \quad (24)$$

有了  $\mathbf{g}_n$  后通过朝前推的 Lanczos 过程, 就可以求出  $B_{1,n}$  的全部元素.

例: 对于 3 个质点和 3 个弹簧的系统, 已知

$$\lambda_1=2, \quad \lambda_2=4, \quad \lambda_3=6,$$

$$\mu_1=3, \quad \mu_2=5,$$

求  $k_1, k_2, k_3, m_2, m_3$  用  $m_1$  的表示式.

$$\text{解} \quad s_{31}^2 = \frac{(\mu_1^2 - \lambda_1^2)(\mu_2^2 - \lambda_1^2)}{(\lambda_2^2 - \lambda_1^2)(\lambda_3^2 - \lambda_1^2)} = 0.2734375,$$

$$s_{32}^2 = \frac{(\mu_1^2 - \lambda_2^2)(\mu_2^2 - \lambda_2^2)}{(\lambda_1^2 - \lambda_2^2)(\lambda_3^2 - \lambda_2^2)} = 0.2625,$$

$$s_{33}^2 = \frac{(\mu_1^2 - \lambda_3^2)(\mu_2^2 - \lambda_3^2)}{(\lambda_1^2 - \lambda_3^2)(\lambda_2^2 - \lambda_3^2)} = 0.4640625,$$

$$\text{得} \quad \mathbf{g}_3 = (0.5229125, 0.5123475, 0.6812213)^T,$$

$$A = \begin{pmatrix} -4 & 0 & 0 \\ 0 & -16 & 0 \\ 0 & 0 & -36 \end{pmatrix},$$

$$\alpha_3 = (A\mathbf{g}_3, \mathbf{g}_3) = -21.99999784,$$

$$\beta_2 = \|A\mathbf{g}_3 - \alpha_3\mathbf{g}_3\| = 13.747726,$$

$$\begin{aligned}
\mathbf{g}_2 &= (\mathbf{A}\mathbf{g}_1 - a_1\mathbf{g}_1)/\beta_2 \\
&= (0.6846531, 0.2236067, -0.6937220)^T, \\
a_2 &= (\mathbf{A}\mathbf{g}_2, \mathbf{g}_2) = -20.0000, \\
\beta_1 &= \|\mathbf{A}\mathbf{g}_2 - a_2\mathbf{g}_2 - \beta_2\mathbf{g}_3\| = 7.4161985, \\
\mathbf{g}_1 &= (\mathbf{A}\mathbf{g}_2 - a_2\mathbf{g}_2 - \beta_2\mathbf{g}_3)/\beta_1 \\
&= (0.5077524, -0.8291561, 0.2338539)^T, \\
a_1 &= (\mathbf{A}\mathbf{g}_1, \mathbf{g}_1) = -14.0000, \\
\text{得 } B_{1,n} &= \begin{pmatrix} -14.0000 & 7.4161985 & 0 \\ 7.4161985 & -20.0000 & 13.747726 \\ 0 & 13.747726 & -22.0000 \end{pmatrix},
\end{aligned}$$

再用递推公式(23)得

$$\begin{aligned}
k_1 &= -a_1m_1 = 14m_1, \\
k_2 &= \beta_1^2m_1^2/k_1 = 3.9285714m_1, \\
m_2 &= -k_2/(a_2 + k_2/m_1) = 0.24444444m_1, \\
k_3 &= \beta_2^2m_2^2/k_2 = 2.8746655m_1, \\
m_3 &= -k_3/(a_3 + k_3/m_2) = 0.280729m_1,
\end{aligned}$$

从公式(24)可以看出, 已知的固有频率必须满足

$$\lambda_1 < \mu_1 < \lambda_2 < \mu_2 < \cdots < \mu_{n-1} < \lambda_n,$$

否则  $s_{n,j}$  中会出现虚数, 问题就没有解.

## § 5 解对称线性代数方程组的 Lanczos 算法

在第1章介绍的共轭斜量法, 适用于系数矩阵是对称正定的或对称非负的方程组的解法, 在实际问题中也要碰到系数矩阵对称不定的方程组, 即系数矩阵有正的特征值, 也有负的特征值, 例如文献[17]中讨论的就是一种对称不定的方程组. 对于对称不定的方程组, 共轭斜量法失效了, 首先因为



$(Ax, y)$  此时不是  $x, y$  的一种内积。

C. C. Paige and M. A. Saunders 在 [18] 和 B. N. Parlett 在 [19] 都介绍了用 Lanczos 过程来解这种方程组的方法, 称为 Lanczos 算法。对于方程组

$$Ax = b, \quad (25)$$

$A$  是  $n \times n$  实对称矩阵, Lanczos 算法的程式如下:

1. 取一个初始向量  $x_0$ , 计算  $r_0 = b - Ax_0$ ;
2. 求  $\|r_0\|$  和  $q_1 = r_0 / \|r_0\|$ ;
3. 由 Lanczos 过程 (7),

$$\beta_i q_{i+1} = Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1},$$

$$i = 1, 2, \dots, j,$$

$$\beta_0 q_0 = 0,$$

求标准正交化序列  $q_1, q_2, \dots, q_{j+1}, \alpha_1, \alpha_2, \dots, \alpha_j, \beta_1, \beta_2, \dots, \beta_j$ , 其中

$$\alpha_i = (Aq_i, q_i),$$

$$\beta_i = \|Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1}\|;$$

4. 由 
$$T_j = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ 0 & & \ddots & \beta_{j-1} & \alpha_j \\ & & & \beta_{j-1} & \alpha_j \end{pmatrix},$$

如果  $\det(T_j) \neq 0$ , 求方程组

$$T_j f_j = \|r_0\| e_1 \quad (26)$$

的解  $f_j = (\delta_1^{(j)}, \delta_2^{(j)}, \dots, \delta_j^{(j)})^T$ ;

5. 记矩阵  $Q_j = [q_1, q_2, \dots, q_j]$ , 构造  $Q_j f_j$ , 则  $x_j = Q_j f_j + x_0$  是方程 (25) 的第  $j$  次近似解。并成立误差估计式

$$\|Ax_j - b\| = |\delta_j^{(j)}| \beta_j, \quad (27)$$

(27)式可以证明如下: 过程(7)成立关系式

$$AQ_j - Q_j T_j = \beta_j q_{j+1} e_j^T,$$

于是

$$AQ_j f_j - Q_j T_j f_j = \beta_j q_{j+1} \delta_j^{(j)},$$

$$A(x_j - x_0) - \|r_0\| q_1 = \beta_j \delta_j^{(j)} q_{j+1},$$

$$Ax_j - b = \beta_j \delta_j^{(j)} q_{j+1},$$

故

$$\|Ax_j - b\| = \beta_j |\delta_j^{(j)}|$$

成立。

上述 Lanczos 方法是将求方程组 (25) 的解的问题, 化成求对称三对角方程组 (26) 的问题。从理论上来说由于  $q_1, q_2, \dots, q_n$  的互相正交性知当  $j=n$  时,  $\beta_n=0$ , 即  $x_n$  为正确解; 但在实际计算上, 由于  $q_1, q_2, \dots, q_n$  的正交性消失,  $\beta_n$  不一定为 0, 但很多计算表明, 常常有某些  $j < n$ , 使  $\beta_j |\delta_j^{(j)}|$  很小, 参见 [18], [19], [13]。因此 Lanczos 方法求解对称方程组, 仍是一种有效的方法, 特别它可以充分利用矩阵  $A$  的稀疏性, 用到  $A$  的时候只要求计算  $Av$ , 即  $A$  与向量  $v$  的乘法, 这一点不但可以节约存贮单元, 而且可以容易被向量运算机来实现。

在计算过程中, 如果预先知道  $\beta_j \delta_j^{(j)}$ , 那么只要等到  $\beta_j \delta_j^{(j)}$  足够小时, 才进行步骤 4 和 5, 就可大大节省计算量, 因此  $\beta_j \delta_j^{(j)}$  记为  $\Delta_j$  称为这一算法的判据。于是怎样不通过解方程组 (26), 而直接计算  $\delta_j^{(j)}$  的问题就提出来了。文献 [18] 中也提出了一种直接计算  $\delta_j^{(j)}$  的方法, 在本章 § 3 定理 2.6 给出一个更简单的求  $\delta_j^{(j)}$  的公式, 即为当  $\alpha = \|r_0\|$  时代入公式 (11) 得

$$\Delta_j = \beta_j \delta_j^{(j)} = (-1)^{j-1} \|r_0\| \frac{\beta_1 \beta_2 \cdots \beta_j}{d_j},$$

并有递推公式

$$\Delta_j = -\Delta_{j-1} \frac{d_{j-1}}{d_j} \beta_j,$$

行列式  $d_j$  可以通过递推公式

$$d_i = \alpha_i d_{i-1} - \beta_{i-1}^2 d_{i-2}, \quad \beta_0 = 0$$

得到. 如果  $d_j = 0$  的话  $d_{j+1}$  不为 0,  $d_{j-1}$  也不为 0, 可以计算  $\Delta_{j+1}$ ,

$$\Delta_{j+1} = \Delta_{j-1} \frac{d_{j-1}}{d_{j+1}} \beta_j \beta_{j+1} = -\Delta_{j-1} \frac{\beta_{j+1}}{\beta_j}.$$

对于  $\|Ax_j - b\|$ , 除了 (27) 所示的等式外, 还可以利用  $\beta_1 \beta_2 \cdots \beta_j$  的极值性质 (定理 2.5) 得到下述估计

$$\|Ax_j - b\| = \frac{\|r_0\|}{|d_j|} \min_{\psi \in mP_j} \|\psi(T_k) e_1\|,$$

$k$  是大于  $j$ 、小于等于  $n$  的任意正整数.

在 [18] 中指出, 当矩阵  $A$  是对称正定时, Lanczos 算法得到的  $x_j$  与共轭斜量法第  $j$  步近似  $x_j$  是相同的. 这一结论, 我们证明如下: 记方程 (25) 的正确解为  $\tilde{x}$ , Lanczos 算法得到的第  $j$  次近似解为  $x_j$ , 只要证明

$$(A(\tilde{x} - x_j), A^l r_0) = 0, \quad l = 0, 1, 2, \dots, j-1 \quad (28)$$

就行了. 因  $A$  正定, 假定 Lanczos 算法前  $j$  步得到的  $\beta_1, \beta_2, \dots, \beta_{j-1}$  都不为 0, 此时  $Q_j = [q_1, q_2, \dots, q_j]$  可以算得,  $T_j = Q_j^T A Q_j$  也是一个正定矩阵, 因此可以将  $T_j$  进行三角分解

$$T_j = L_j D_j L_j^T,$$

其中  $L_j$  是下三角阵, 对角元为 1,  $D_j$  为对角阵. 记

$$P_j = Q_j L_j^{-T} = [p_1, p_2, \dots, p_j],$$

$p_i$  是  $q_1, q_2, \dots, q_i$  的线性组合, 而  $q_i$  是  $r_0, Ar_0, \dots, A^{i-1}r_0$  的线性组合, 因此  $p_i$  是  $r_0, Ar_0, \dots, A^{i-1}r_0$  的线性组合. 反之由  $\beta_1, \beta_2, \dots, \beta_{j-1}$  不为 0, 可知  $A^{l-1}r_0$  是  $q_1, q_2, \dots, q_l$  的

线性组合,从而也是  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_l$  的线性组合. 这样(28)式等价于

$$P_j^T A(\tilde{\mathbf{x}} - \mathbf{x}_j) = 0, \quad (29)$$

因为

$$\mathbf{x}_j = Q_j \mathbf{f}_j + \mathbf{x}_0,$$

故

$$A(\tilde{\mathbf{x}} - \mathbf{x}_j) = \mathbf{r}_0 - A Q_j \mathbf{f}_j,$$

所以

$$\begin{aligned} P_j^T A(\tilde{\mathbf{x}} - \mathbf{x}_j) &= P_j^T \mathbf{r}_0 - P_j^T A Q_j \mathbf{f}_j \\ &= L_j^{-1} Q_j^T \mathbf{r}_0 - L_j^{-1} Q_j^T A Q_j \mathbf{f}_j \\ &= L_j^{-1} \|\mathbf{r}_0\| \mathbf{e}_1 - L_j^{-1} T_j \mathbf{f}_j = 0, \end{aligned}$$

这就证明了  $\mathbf{x}_j$  也是共轭斜量法得到的第  $j$  步近似解. 如果  $\beta_l = 0$ , 那么由(27)知  $\mathbf{x}_l$  即为  $\tilde{\mathbf{x}}$ , 此时  $\mathbf{r}_0, A\mathbf{r}_0, \dots, A^l \mathbf{r}_0$  线性相关, 从而共轭斜量法得到的第  $l$  步近似也是  $\tilde{\mathbf{x}}$ . 这就证明了, 当  $A$  正定时, 两种方法是等价的.

$$\text{另外由 } P_j^T A P_j = L_j^{-1} Q_j^T A Q_j L_j^{-T} = D_j,$$

也即  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_j$  是  $A$ -正交的向量组, 容易知道共轭斜量法中得到的向量组  $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$  中的  $\mathbf{q}_i$  与  $\mathbf{p}_i$  至多差个常数倍.

在矩阵  $A$  对称不定的情况, Lanczos 算法解方程组, 还有两个问题需要考虑. 第一个问题是此时  $T_j$  不能保证可以进行三角分解,  $T_j = L_j D_j L_j^T$ , 因为  $T_j$  不是正定, 它的某些顺序主子式可能为 0, 因此如何对方程组(26)求解就是一个问题了. 当然可以采用选主元的消去法, 不过那样一来,  $T_j$  的三对角的形状会遭破坏, 存储量就要增加.

第二个问题是第  $j$  次近似  $\mathbf{x}_j$  与第  $j+1$  次近似  $\mathbf{x}_{j+1}$ , 不象共轭斜量法中那样有简单的表示式  $\mathbf{x}_{j+1} = \mathbf{x}_j + \alpha_j \mathbf{q}_j$ , 这是因为  $\mathbf{f}_{j+1}$  的首  $j$  个分量与  $\mathbf{f}_j$  没有简单的关系, 因此当已经

计算得  $x_j$  后, 如果还要计算  $x_{j+1}$ , 就不能充分利用  $x_j$  提供的信息, 这也是很大的浪费. 怎样避免这种浪费, 即充分利用已算得的  $x_j$ , 也是一个问题.

O. C. Paige 和 A. Saunders 在 [18] 中提出了一种称为 SYMMLQ 的方法, 就是解决上述两个问题的一种方法. 下面介绍 SYMMLQ 方法:

在 Lanczos 算法步骤 1、2、3 以后, 对每个  $j$  得到对称三对角矩阵  $T_j$ ; 将  $T_j$  进行 LQ 分解

$$T_j = \tilde{L}_j P_j,$$

其中  $L_j$  是下三角矩阵,  $P_j$  是正交矩阵. 为了求  $\tilde{L}_j$  和  $P_j$ , 对  $T_j$  右乘一系列旋转阵  $H_{i,i+1}$ , 使得

$$T_j H_{12} H_{23} \cdots H_{j-1,j} = \tilde{L}_j,$$

其中  $H_{i,i+1} = (h_{ik})$ , 除  $h_{ii} = c_i$ ,  $h_{i+1,i+1} = -c_i$ ,  $h_{i,i+1} = h_{i+1,i} = s_i$  外其余元素都与单位阵的相同. 右乘  $H_{i,i+1}$  的目的是为了将  $T_j$  的第 1 条超对角线上第  $i$  个元素化成 0. 我们以 5 阶矩阵为例来观察计算的程式.

$$T_5 = \begin{pmatrix} \alpha_1 & \beta_1 & 0 & 0 & 0 \\ \beta_1 & \alpha_2 & \beta_2 & 0 & 0 \\ 0 & \beta_2 & \alpha_3 & \beta_3 & 0 \\ 0 & 0 & \beta_3 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix},$$

$$T_5 H_{12} = \begin{pmatrix} \alpha_1 & \beta_1 & 0 & 0 & 0 \\ \beta_1 & \alpha_2 & \beta_2 & 0 & 0 \\ 0 & \beta_2 & \alpha_3 & \beta_3 & 0 \\ 0 & 0 & \beta_3 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix} \begin{pmatrix} c_1 & s_1 & 0 & 0 & 0 \\ s_1 & -c_1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$= \begin{pmatrix} \alpha_1 c_1 + \beta_1 s_1 & \alpha_1 s_1 - \beta_1 c_1 & 0 & 0 & 0 \\ \beta_1 c_1 + \alpha_2 s_1 & \beta_1 s_1 - \alpha_2 c_1 & \beta_2 & 0 & 0 \\ \beta_2 s_1 & -\beta_2 c_1 & \alpha_3 & \beta_3 & 0 \\ 0 & 0 & \beta_3 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix},$$

于是取  $s_1 = \frac{\beta_1}{\sqrt{\alpha_1^2 + \beta_1^2}}, c_1 = \frac{\alpha_1}{\sqrt{\alpha_1^2 + \beta_1^2}},$

就能把  $\beta_1$  位置上元素化成 0, 记

$$\gamma_1 = \alpha_1 c_1 + \beta_1 s_1, \quad \delta_2 = \beta_1 c_1 + \alpha_2 s_1, \quad \varepsilon_3 = \beta_2 s_1,$$

$$\tilde{\gamma}_2 = \beta_1 s_1 - \alpha_2 c_1, \quad \tilde{\delta}_3 = -\beta_2 c_1,$$

于是  $T_5 H_{12} = \begin{pmatrix} \gamma_1 & & & & \\ \delta_2 & \tilde{\gamma}_2 & \beta_2 & & \\ \varepsilon_3 & \tilde{\delta}_3 & \alpha_3 & \beta_3 & \\ & & \beta_3 & \alpha_4 & \beta_4 \\ & & & \beta_4 & \alpha_5 \end{pmatrix},$

$$T_5 H_{12} H_{23} = \begin{pmatrix} \gamma_1 & & & & \\ \delta_2 & \tilde{\gamma}_2 & \beta_2 & & \\ \varepsilon_3 & \tilde{\delta}_3 & \alpha_3 & \beta_3 & \\ & & \beta_3 & \alpha_4 & \beta_4 \\ & & & \beta_4 & \alpha_5 \end{pmatrix} \begin{pmatrix} 1 & & & & \\ & c_2 & s_2 & & \\ & s_2 & -c_2 & & \\ & & & 1 & \\ & & & & 1 \end{pmatrix}$$

$$= \begin{pmatrix} \gamma_1 & 0 & & 0 & 0 \\ \delta_2 & \tilde{\gamma}_2 c_2 + \beta_2 s_2 & \tilde{\gamma}_2 s_2 - \beta_2 c_2 & 0 & 0 \\ \varepsilon_3 & \tilde{\delta}_3 c_2 + \alpha_3 s_2 & \tilde{\delta}_3 s_2 - \alpha_3 c_2 & \beta_3 & 0 \\ 0 & \beta_3 s_2 & -\beta_3 c_2 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix},$$

取  $s_2 = \beta_2 / (\tilde{\gamma}_2^2 + \beta_2^2)^{\frac{1}{2}}, c_2 = \tilde{\gamma}_2 / (\tilde{\gamma}_2^2 + \beta_2^2)^{\frac{1}{2}},$

则  $\beta_2$  位置上的元素化成 0. 记

$$\gamma_2 = \tilde{\gamma}_2 c_2 + \beta_2 s_2,$$

$$\delta_2 = \tilde{\delta}_2 c_2 + \alpha_2 s_2,$$

$$\varepsilon_2 = \beta_2 s_2,$$

$$\tilde{\gamma}_3 = \tilde{\delta}_3 s_2 - \alpha_3 c_2,$$

$$\tilde{\delta}_4 = -\beta_3 c_2,$$

从上述计算中可见:  $\gamma$ 、 $\delta$ 、 $\varepsilon$  在下次旋转时不改变, 而  $\tilde{\gamma}$ 、 $\tilde{\delta}$  在下次旋转时要改变。

$$T_5 H_{12} H_{23} = \begin{pmatrix} \gamma_1 & 0 & 0 & 0 & 0 \\ \delta_2 & \gamma_2 & 0 & 0 & 0 \\ \varepsilon_3 & \delta_3 & \tilde{\gamma}_3 & \beta_3 & 0 \\ 0 & \varepsilon_4 & \tilde{\delta}_4 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix},$$

$$T_5 H_{12} H_{23} H_{34} = \begin{pmatrix} \gamma_1 & 0 & 0 & 0 & 0 \\ \delta_2 & \gamma_2 & 0 & 0 & 0 \\ \varepsilon_3 & \delta_3 & \tilde{\gamma}_3 & \beta_3 & 0 \\ 0 & \varepsilon_4 & \tilde{\delta}_4 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{pmatrix} \begin{pmatrix} 1 & & & & \\ & 1 & & & \\ & & c_3 & s_3 & \\ & & s_3 & -c_3 & \\ & & & & 1 \end{pmatrix}$$

$$= \begin{pmatrix} \gamma_1 & 0 & 0 & 0 & 0 \\ \delta_2 & \gamma_2 & 0 & 0 & 0 \\ \varepsilon_3 & \delta_3 & \tilde{\gamma}_3 c_3 + \beta_3 s_3 & \tilde{\gamma}_3 s_3 - \beta_3 c_3 & 0 \\ 0 & \varepsilon_4 & \tilde{\delta}_4 c_3 + \alpha_4 s_3 & \tilde{\delta}_4 s_3 - \alpha_4 c_3 & \beta_4 \\ 0 & 0 & \beta_4 s_3 & -\beta_4 c_3 & \alpha_5 \end{pmatrix},$$

要使 (3, 4) 位置上的元素为 0, 取

$$s_3 = \beta_3 / (\tilde{\gamma}_3^2 + \beta_3^2)^{\frac{1}{2}}, \quad c_3 = \tilde{\gamma}_3 / (\tilde{\gamma}_3^2 + \beta_3^2)^{\frac{1}{2}},$$

记

$$\gamma_3 = \tilde{\gamma}_3 c_3 + \beta_3 s_3,$$

$$\delta_4 = \tilde{\delta}_4 c_3 + \alpha_4 s_3,$$

$$\varepsilon_5 = \beta_4 s_3,$$

$$\begin{aligned}\tilde{\gamma}_4 &= \tilde{\delta}_4 s_3 - \alpha_4 c_3, \\ \tilde{\delta}_5 &= -\beta_4 c_3, \\ T_5 H_{12} H_{23} H_{34} &= \begin{pmatrix} \gamma_1 & 0 & 0 & 0 & 0 \\ \delta_2 & \gamma_2 & 0 & 0 & 0 \\ \varepsilon_3 & \delta_3 & \gamma_3 & 0 & 0 \\ 0 & \varepsilon_4 & \delta_4 & \tilde{\gamma}_4 & \beta_4 \\ 0 & 0 & \varepsilon_5 & \tilde{\delta}_5 & \alpha_5 \end{pmatrix},\end{aligned}$$

总结上述过程, 有一般的计算程式

$$\tilde{\gamma}_1 = \alpha_1, \quad \tilde{\delta}_2 = \beta_1,$$

则使用  $H_{k,k+1}$  旋转时

$$s_k = \beta_k / (\tilde{\gamma}_k^2 + \beta_k^2)^{1/2}, \quad c_k = \tilde{\gamma}_k / (\tilde{\gamma}_k^2 + \beta_k^2)^{1/2},$$

通过旋转得到

$$\left. \begin{aligned}\gamma_k &= \tilde{\gamma}_k c_k + \beta_k s_k = (\tilde{\gamma}_k^2 + \beta_k^2)^{1/2}, \\ \delta_{k+1} &= \tilde{\delta}_{k+1} c_k + \alpha_{k+1} s_k, \\ \varepsilon_{k+2} &= \beta_{k+1} s_k, \\ \tilde{\gamma}_{k+1} &= \tilde{\delta}_{k+1} s_k - \alpha_{k+1} c_k, \\ \tilde{\delta}_{k+2} &= -\beta_{k+1} c_k, \\ k &= 1, 2, \dots,\end{aligned}\right\} \quad (30)$$

按此公式计算

$$T_5 H_{12} H_{23} H_{34} H_{45} = \begin{pmatrix} \gamma_1 & & & & \\ \delta_2 & \gamma_2 & & & \\ \varepsilon_3 & \delta_3 & \gamma_3 & & \\ 0 & \varepsilon_4 & \delta_4 & \gamma_4 & \\ 0 & 0 & \varepsilon_5 & \delta_5 & \tilde{\gamma}_5 \end{pmatrix} = \tilde{L}_5,$$

而

$$P_5 = H_{45} H_{34} H_{23} H_{12}.$$

如果  $\beta_k \neq 0$ , 则有  $\gamma_k > 0$ , 但  $\tilde{\gamma}_{k+1}$  可能为 0.

从上述推导还可知



$$T_6 P_5^T = T_6 H_{12} H_{23} H_{34} H_{45} = \begin{pmatrix} \gamma_1 & & & & & \\ \delta_2 & \gamma_2 & & & & \\ \varepsilon_3 & \delta_3 & \gamma_3 & & & \\ 0 & \varepsilon_4 & \delta_4 & \gamma_4 & & \\ 0 & 0 & \varepsilon_5 & \delta_5 & \tilde{\gamma}_5 & \beta_5 \\ 0 & 0 & 0 & \varepsilon_6 & \tilde{\delta}_6 & \alpha_6 \end{pmatrix},$$

其中  $\tilde{\delta}_6 = -\beta_5 \varepsilon_4, \quad \varepsilon_6 = \beta_5 \delta_4,$

而  $P_6^T = P_5^T H_{56}.$

$T_j$  经过  $LQ$  分解后,  $T_j = \tilde{L}_j P_j$ , 于是方程组 (26)

$$T_j f_j = \|r_0\| e_1$$

可化成  $\tilde{L}_j P_j f_j = \|r_0\| e_1,$

令  $\tilde{z}_j = P_j f_j$ , 如果  $\tilde{\gamma}_j \neq 0$ , 即  $\det(\tilde{L}_j) \neq 0$ , 那么可以由方程

$$\tilde{L}_j \tilde{z}_j = \|r_0\| e_1 \quad (31)$$

唯一确定  $\tilde{z}_j$ . 如果知道了  $\beta_j$ , 由  $\gamma_j = (\tilde{\gamma}_j^2 + \beta_j^2)^{1/2}$ , 可构造一个新的下三角阵

$$L_j = \begin{pmatrix} \gamma_1 & & & & & \\ \delta_2 & \gamma_2 & & & & \\ \varepsilon_3 & \delta_3 & \ddots & & & \\ & \varepsilon_4 & \ddots & \ddots & & \\ & & \ddots & \ddots & \gamma_{j-1} & \\ & & & \varepsilon_j & \delta_j & \gamma_j \end{pmatrix},$$

与  $\tilde{L}_j$  只差一个对角元, 当  $\beta_j \neq 0$  时  $\det(L_j) \neq 0$ , 如果令  $z_j = (\zeta_1, \zeta_2, \dots, \zeta_j)^T$ ,  $z_j$  可由

$$L_j z_j = \|r_0\| e_1$$

唯一确定. 因为  $L_j$  是下三角阵, 因此由

$$L_{j+1} z_{j+1} = \|r_0\| e_1$$

确定的  $z_{j+1}$  必有  $z_{j+1} = (z_j^T, \zeta_{j+1})^T$ , 即  $z_{j+1}$  的为首  $j$  个分量

为  $z_j$ 。同时由方程(31), 可知

$$\tilde{z}_j = (\zeta_1, \zeta_2, \dots, \zeta_{j-1}, \tilde{\zeta}_j)^T,$$

并且当  $\tilde{\gamma}_j \neq 0$  时有

$$\tilde{\gamma}_j \tilde{\zeta}_j = \gamma_j \zeta_j.$$

再由  $x_j = Q_j f_j + x_0 = Q_j P_j^T \tilde{z}_j + x_0$ ,

需要分析  $Q_j P_j^T$ 。由

$$\begin{aligned} Q_{j+1} P_{j+1}^T &= (Q_j, q_{j+1}) \begin{pmatrix} P_j^T & 0 \\ 0 & 1 \end{pmatrix} H_{j,j+1} \\ &= (Q_j P_j^T, q_{j+1}) \begin{pmatrix} I & 0 \\ 0 & c_j & s_j \\ 0 & s_j & -c_j \end{pmatrix}, \end{aligned} \quad (32)$$

可知  $Q_{j+1} P_{j+1}^T$  的首  $j-1$  列与  $Q_j P_j^T$  为首  $j-1$  列是完全相同的。记

$$Q_j P_j^T = (w_1, w_2, \dots, w_{j-1}, \tilde{w}_j),$$

于是  $Q_{j+1} P_{j+1}^T = (w_1, w_2, \dots, w_{j-1}, w_j, \tilde{w}_{j+1})$ ,

比较(32)的两边有

$$\left. \begin{aligned} w_j &= c_j \tilde{w}_j + s_j q_{j+1}, \\ \tilde{w}_{j+1} &= s_j \tilde{w}_j - c_j q_{j+1}, \end{aligned} \right\} \quad (33)$$

因为  $P_1^T = 1$ , 因此  $\tilde{w}_1 = q_1$ 。

如记  $x_j^L = (w_1, w_2, \dots, w_j) z_j + x_0$ ,

就有

$$x_j^L = x_{j-1}^L + \zeta_j w_j, \quad (34)$$

而 Lanczos 算法的近似解

$$x_j = (w_1, w_2, \dots, w_{j-1}, \tilde{w}_j) \tilde{z}_j + x_0 = x_{j-1}^L + \tilde{\zeta}_j \tilde{w}_j. \quad (35)$$

这就是 SYMMLQ 方法求  $x_j$  的过程, 在实际计算中,  $x_j^L$  也是一个近似解, 在迭代过程中, 可以在前面只用(34),

求  $\mathbf{x}_1^L, \mathbf{x}_2^L, \dots, \mathbf{x}_{i-1}^L$ , 最后一步才用 (35) 求  $\mathbf{x}_i$ , 归纳一下 SYMMLQ 方法的步骤如下:

1. 取初始近似向量  $\mathbf{x}_0$ , 求  $\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0$ ,

计算  $\|\mathbf{r}_0\|$ ,  $\mathbf{q}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|$ ,  $\tilde{\mathbf{w}}_1 = \mathbf{q}_1$ ,

$\alpha_1 = (A\mathbf{q}_1, \mathbf{q}_1)$ ,  $\tilde{\gamma}_1 = \alpha_1$ , 求  $A\mathbf{q}_1 - \alpha_1\mathbf{q}_1$ ,

计算  $\|A\mathbf{q}_1 - \alpha_1\mathbf{q}_1\| = \beta_1$ ,

如果  $\beta_1 = 0$ , 则  $\mathbf{x}_0 + \frac{\|\mathbf{r}_0\|}{\alpha_1} \mathbf{q}_1 \rightarrow \mathbf{x}$ , 转 5,

否则计算  $\mathbf{q}_2 = (A\mathbf{q}_1 - \alpha_1\mathbf{q}_1) / \beta_1$ ,  $(A\mathbf{q}_2, \mathbf{q}_2) = \alpha_2$ ,

$d_1 = \alpha_1$ , 如果  $d_1 \neq 0$ , 则计算  $\Delta_1 = \frac{\|\mathbf{r}_0\| \beta_1}{|d_1|}$ ,

$d_2 = \alpha_1\alpha_2 - \beta_1^2$ ,  $\varepsilon_1 = \varepsilon_2 = \delta_1 = 0$ ,  $\tilde{\delta}_2 = \beta_1$ ,

$\gamma_1 = (\tilde{\gamma}_1^2 + \beta_1^2)^{1/2}$ ,  $\zeta_1 = \|\mathbf{r}_0\| / \gamma_1$ ,

$s_1 = \beta_1 / \gamma_1$ ,  $c_1 = \tilde{\gamma}_1 / \gamma_1$ ,

$\tilde{\gamma}_2 = \tilde{\delta}_2 s_1 - \alpha_2 c_1$ ,

$\mathbf{w}_1 = c_1 \tilde{\mathbf{w}}_1 + s_1 \mathbf{q}_2$ ,  $\tilde{\mathbf{w}}_2 = s_1 \tilde{\mathbf{w}}_1 - c_1 \mathbf{q}_2$ ,

$\mathbf{x}_1^L = \mathbf{x}_0 + \zeta_1 \mathbf{w}_1$ ,

求  $A\mathbf{q}_2 - \alpha_2\mathbf{q}_2 - \beta_1\mathbf{q}_1$ ,

计算  $\beta_2 = \|A\mathbf{q}_2 - \alpha_2\mathbf{q}_2 - \beta_1\mathbf{q}_1\|$ ,

如果  $d_2 \neq 0$  及  $\beta_2 = 0$ , 则  $\tilde{\zeta}_2 = -\delta_2 \zeta_1 / \tilde{\gamma}_2$ ,

$\mathbf{x}_2 = \mathbf{x}_1^L + \tilde{\zeta}_2 \tilde{\mathbf{w}}_2 \rightarrow \mathbf{x}$  转 5.

如果  $\beta_2 \neq 0$ , 计算  $\Delta_2 = \frac{\|\mathbf{r}_0\| \beta_1 \beta_2}{|d_2|}$

$2 \rightarrow i$  转 2.

如果  $d_2 = 0$ ,  $2 \rightarrow i$  转 2.

2.  $\varepsilon_{i+1} = \beta_i s_{i-1}$ ,  $\tilde{\delta}_{i+1} = -\beta_i c_{i-1}$ ,

$\gamma_i = (\tilde{\gamma}_i^2 + \beta_i^2)^{1/2}$ ,  $s_i = \beta_i / \gamma_i$ ,  $c_i = \tilde{\gamma}_i / \gamma_i$ ,

$\mathbf{q}_{i+1} = (A\mathbf{q}_i - \alpha_i\mathbf{q}_i - \beta_{i-1}\mathbf{q}_{i-1}) / \beta_i$ ,

$$\alpha_{i+1} = (Aq_{i+1}, q_{i+1}),$$

$$\delta_{i+1} = \tilde{\delta}_{i+1}c_i + \alpha_{i+1}s_i,$$

$$\tilde{\gamma}_{i+1} = \tilde{\delta}_{i+1}s_i - \alpha_{i+1}c_i,$$

$$w_i = \tilde{w}_i c_i + q_{i+1} s_i,$$

$$\tilde{w}_{i+1} = \tilde{w}_i s_i - q_{i+1} c_i,$$

$$\zeta_i = \frac{-\varepsilon_i \zeta_{i-2} - \delta_i \zeta_{i-1}}{\gamma_i},$$

$$x_i^L = x_{i-1}^L + \zeta_i w_i.$$

3. 构造  $g_{i+1} = Aq_{i+1} - \alpha_{i+1}q_{i+1} - \beta_i q_i$ ,

计算  $\beta_{i+1} = \|g_{i+1}\|$ ,

计算  $d_{i+1} = \alpha_{i+1}d_i - \beta_i^2 d_{i-1}$ ,

如果  $d_{i+1} = 0$ , 则  $i+1 \rightarrow i$  转 2,

如果  $d_{i+1} \neq 0$ , 计算

$$\Delta_{i+1} = \begin{cases} \Delta_i \beta_{i+1} \left| \frac{d_i}{d_{i+1}} \right|, & \text{当 } d_i \neq 0, \\ \Delta_{i-1} \beta_i \beta_{i+1} \left| \frac{d_{i-1}}{d_{i+1}} \right|, & \text{当 } d_i = 0, \end{cases}$$

如果  $\Delta_{i+1} < \varepsilon$  转 4, 否则  $i+1 \rightarrow i$  转 2.

4. 计算

$$\tilde{\zeta}_{i+1} = -(\varepsilon_{i+1}\zeta_{i-1} + \delta_{i+1}\zeta_i)/\tilde{\gamma}_{i+1},$$

$$x_{i+1} = x_i^L + \tilde{\zeta}_{i+1} \tilde{w}_{i+1} \rightarrow x.$$

5. 输出计算结果  $x$ ; 停止.

## 第 3 章

### 解特征值问题的 QL 方法

计算中小型对称矩阵的特征值的一个比较有效的方法,是先将这个矩阵,通过第 2 章 § 2 所介绍的正交相似变换化成一个对称三对角矩阵,然后再用带特别选取的位移的 QL 方法(或 QR 方法)求得这个对称三对角矩阵的特征值。

熟知的 Givens 方法,是将矩阵化为对称三对角矩阵后,再用对分法求后者的特征值,但是对分法的收敛速度是一次的,而带特别选取的位移的 QL 方法(或 QR 方法),收敛速度一般是三次的,并且每次计算量也是相当于对分法,因此这种 QL 方法是更有效。

本章要介绍 QL 方法有关的理论分析,指出它的收敛速度为何一般是三次。

#### § 1 QL 方法的一般性质

设  $A$  是  $n \times n$  实对称矩阵,  $\sigma$  是实参数,将  $A - \sigma I$  分解成

$$A - \sigma I = QL,$$

其中  $Q$  是  $n \times n$  正交矩阵,  $L$  是下三角阵。记

$$\hat{A} = LQ + \sigma I,$$

于是  $Q\hat{A}Q^* = QL + \sigma I = A - \sigma I + \sigma I = A,$

故

$$\hat{A} = Q^*AQ, \quad (1)$$

可知  $\hat{A}$  仍为实对称矩阵,并且正交相似于  $A$ 。

给定实数  $\sigma_k, k=1, 2, \dots$ , 令  $A_1=A$ ,

$$A_k - \sigma_k I = Q_k L_k, \quad (2)$$

$Q_k$  是正交阵  $L_k$  是下三角阵,

$$A_{k+1} = L_k Q_k + \sigma_k I \quad (3)$$

形成了矩阵序列  $\{A_k\}$ , 如果  $A_k$  趋于一个对角阵, 那么  $A$  的全部特征值就得到了. 这样求  $A$  的特征值的方法, 就称为带位移  $\{\sigma_k\}$  的 QL 方法. 如果在 (2) 中的分解是

$$A_k - \sigma_k I = Q_k R_k,$$

$Q_k$  是正交阵,  $R_k$  是上三角阵, 且

$$A_{k+1} = R_k Q_k + \sigma_k I,$$

那么这样的求特征值的方法, 称为带位移  $\{\sigma_k\}$  的 QR 方法.

现在来看 QL 分解是怎么进行的, 记  $Q = [q_1, q_2, \dots, q_n]$ ,  $A - \sigma I = [a_1, a_2, \dots, a_n]$ ,

$$L = \begin{pmatrix} l_{11} & & & 0 \\ l_{21} & l_{22} & & \\ l_{31} & l_{32} & l_{33} & \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & l_{n3} & \dots & l_{nn} \end{pmatrix},$$

比较  $A - \sigma I = QL$  两边知

$$a_n = l_{nn} q_n,$$

因  $\|q_n\|=1$ , 故  $l_{nn} = \|a_n\|$ , 如果取  $l_{nn} > 0$ , 那么由  $a_n$  唯一确定  $l_{nn}$  和  $q_n$ . 再由

$$a_{n-1} = q_n l_{n, n-1} + q_{n-1} l_{n-1, n-1},$$

或

$$q_{n-1} l_{n-1, n-1} = a_{n-1} - q_n l_{n, n-1},$$

根据  $(q_{n-1}, q_n) = 0$ , 得

$$l_{n, n-1} = (a_{n-1}, q_n),$$

同样确定

$$l_{n-1, n-1} = \|a_{n-1} - q_n l_{n, n-1}\|,$$

$$q_{n-1} = (a_{n-1} - q_n l_{n,n-1}) / l_{n-1,n-1},$$

一般地有

$$q_k l_{kk} = a_k - q_{k+1} l_{k+1,k} - q_{k+2} l_{k+2,k} - \cdots - q_n l_{n,k},$$

其中

$$l_{jk} = (a_k, q_j),$$

$$j = k+1, k+2, \dots, n,$$

$$l_{k,k} = \left\| a_k - \sum_{j=k+1}^n q_j l_{jk} \right\|,$$

$$q_k = \left( a_k - \sum_{j=k+1}^n q_j l_{jk} \right) / l_{k,k},$$

$$k = n-1, n-2, \dots, 1.$$

如果  $l_{k,k} = 0$ , 那么上述过程就要中断. 这时

$$a_k - \sum_{j=k+1}^n q_j l_{jk} = 0,$$

也即  $a_k, a_{k+1}, \dots, a_n$  线性相关, 于是  $\det(A - \sigma I) = 0$ , 说明  $\sigma$  是  $A$  的一个特征值.

反之如果  $\sigma$  是  $A$  的一个特征值, 则至少存在一个  $k$ , 使  $l_{k,k} = 0$ . 实际上若  $\sigma$  是  $A$  的特征值, 那么  $\det(A - \sigma I) = 0$ , 从而至少有一个  $k$ , 使得  $a_k, a_{k+1}, \dots, a_n$  线性相关. 但

$$l_{kk} q_k = a_k - \sum_{j=k+1}^n q_j l_{jk},$$

$$\text{有} \quad \left( a_k - \sum_{j=k+1}^n q_j l_{jk}, q_i \right) = 0, \quad i = k+1, \dots, n,$$

因此  $\sum_{j=k+1}^n q_j l_{jk}$  是  $\{a_{k+1}, a_{k+2}, \dots, a_n\}$  中对  $a_k$  的最佳逼近,

因此  $a_k - \sum_{j=k+1}^n q_j l_{jk} = 0$ , 从而推出

$$l_{k,k} = 0,$$

于是关于 QL 方法有

**定理 3.1** 存在一个正整数  $k$ , 使  $l_{kk} = 0$  的充分必要条件

是:  $\sigma$  为  $A$  的特征值.

如果  $\sigma$  不是  $A$  的特征值, 那么  $l_{11}, l_{22}, \dots, l_{nn}$  都不为 0, 这时  $A - \sigma I$  的 QL 分解中, 求  $q_n, q_{n-1}, \dots, q_1$  的过程, 即为将  $a_n, a_{n-1}, \dots, a_1$  正交化的过程. 并且  $Q$  和  $L$  是唯一确定的.

在计算中如果遇到  $l_{kk}=0$ , 但计算还要进行下去, 那么可取  $q_k$  为任意与  $q_{k+1}, q_{k+2}, \dots, q_n$  正交的单位向量. 然后继续求  $q_{k-1}, \dots, q_1$ .

QL 方法的第二个重要性质是:

**定理 3.2** 如果  $A$  是带形矩阵, 则  $\hat{A}$  也是同样带形的矩阵, 即保持原来的带宽.

**证明** 由

$$\begin{aligned} q_n &= \frac{1}{l_{nn}} a_n, \\ q_{n-1} &= \frac{1}{l_{n-1, n-1}} (a_{n-1} - l_{n, n-1} q_n), \\ &\dots\dots\dots \\ q_k &= \frac{1}{l_{kk}} \left( a_k - \sum_{j=k+1}^n q_j l_{jk} \right), \end{aligned}$$

可知  $Q = [q_1, q_2, \dots, q_n]$  的对角线以上部分是保持与  $A$  相同的带宽. 因为  $L$  是下三角阵, 因此  $LQ$  在对角线以上部分仍保持与  $A$  相同的带宽. 但  $\hat{A} = LQ + \sigma I$  是一个对称矩阵, 因此它的对角线以下部分, 也保持与  $A$  相同的带宽. 证毕.

**推论** 若  $A$  是对称三对角矩阵, 则  $\hat{A}$  也是对称三对角矩阵.

下面考察 QL 方法跟 QR 方法的关系. 对于矩阵  $A - \sigma I$  进行 QR 分解, 即

$$A - \sigma I = QR,$$



$Q$  是正交阵,  $R$  是上三角阵. 对同一个矩阵  $A - \sigma I$ , 进行 QR 分解得到的正交阵  $Q$ , 与进行 QL 分解得到的正交阵显然是不相同的. 因此我们为了区别这两个正交阵, 我们用  $Q_R(B)$  表示对  $B$  进行 QR 分解得到的正交阵,  $R_B$  表此分解的上三角阵, 而用  $Q_L(B)$  表示对  $B$  进行 QL 分解而得到的正交阵,  $L_B$  表此分解的下三角阵.

记  $\tilde{I} = (e_n, e_{n-1}, \dots, e_1)$ , 有  $\tilde{I}^{-1} = \tilde{I}$ .

**引理 3.1** 设  $\det(B) \neq 0$ ,  $\tilde{B} = \tilde{I}B\tilde{I}$ , 则

$$Q_L(\tilde{B}) = \tilde{I}Q_R(B)\tilde{I}, \quad L_{\tilde{B}} = \tilde{I}R_B\tilde{I}.$$

**证明**  $\tilde{B} = \tilde{I}B\tilde{I} = \tilde{I}Q_R(B)R_B\tilde{I} = \tilde{I}Q_R(B)\tilde{I}\tilde{I}R_B\tilde{I}$ ,

因为  $\tilde{I}R_B\tilde{I}$  是一个下三角阵, 而  $\tilde{I}Q_R(B)\tilde{I}$  是一个正交阵, 故上式表示了  $\tilde{B}$  的 QL 分解, 但因  $\det(\tilde{B}) = \det(B) \neq 0$ , 因此 QL 分解是唯一的, 故

$$Q_L(\tilde{B}) = \tilde{I}Q_R(B)\tilde{I}.$$

同时还知道  $L_{\tilde{B}} = \tilde{I}R_B\tilde{I}$ . 证毕.

**定理 3.3** 设矩阵序列  $\{A_k\}$ ,  $\{\tilde{A}_k\}$  分别是使用 QL 方法和 QR 方法所得的序列, 所带的位移都是  $\{\sigma_k\}$ . 如果  $\tilde{A}_1 = \tilde{I}A_1\tilde{I}$ , 则

$$\tilde{A}_k = \tilde{I}A_k\tilde{I}$$

对所有  $k > 1$  成立.

**证明**  $A_1 - \sigma_1 I = Q_L(A_1)L_{A_1}$ ,

$$\tilde{A}_1 - \sigma_1 I = Q_R(\tilde{A}_1)R_{\tilde{A}_1},$$

$$A_2 = L_{A_1}Q_L(A_1) + \sigma_1 I,$$

$$\tilde{I}A_2\tilde{I} = \tilde{I}L_{A_1}\tilde{I}\tilde{I}Q_L(A_1)\tilde{I} + \sigma_1 I$$

$$= R_{\tilde{A}_1}Q_R(\tilde{A}_1) + \sigma_1 I = \tilde{A}_2,$$

同理可证  $k=3, 4, \dots$  有

$$\tilde{A}_k = \tilde{I}A_k\tilde{I}. \text{ 证毕.}$$

从定理 3.3 可知, 对某个矩阵进行 QL 方法求特征值与对这个矩阵进行  $\begin{pmatrix} 1 & 2 & 3 & \cdots & n \\ n & n-1 & n-2 & \cdots & 1 \end{pmatrix}$  置换相似变换后的矩阵进行 QR 方法是一样. 在历史上是先提出 QR 方法的, 但是对于有些矩阵它的元素的数量级小的在左上角, 逐渐向右下角增大, 对于这种矩阵, 计算经验表明, 使用 QL 方法, 常更加有利. 这种矩阵在工程技术问题中, 常常遇到, 因此 QL 方法有独立讨论的必要.

现来看 QL 方法与乘幂法、反迭代法的关系.

对于带位移  $\{\sigma_k\}$  的 QL 方法:

$$A_1 = A,$$

$$A_{k+1} = Q_k^* A_k Q_k = Q_k^* Q_{k-1}^* \cdots Q_1^* A_1 Q_1 Q_2 \cdots Q_k,$$

记

$$P_k = Q_1 Q_2 \cdots Q_k,$$

故

$$A_{k+1} = P_k^* A_1 P_k.$$

乘幂法: 取单位向量  $v_1$  作为初始向量,

$$v_{k+1} = (A - \sigma_k) v_k / v_k, \quad (4)$$

取  $v_k$  使得  $\|v_{k+1}\| = 1$ .

反迭代法: 取单位向量  $u_1$  作为初始向量,

$$(A - \sigma_k) u_{k+1} = u_k \tau_k, \quad (5)$$

取  $\tau_k$  使  $\|u_{k+1}\| = 1$ .

**定理 3.4** 如果序列  $\{\sigma_k\}$  中没有一个是  $A$  的特征值, 并且若  $v_1 = e_n, u_1 = e_1$ , 则

$$v_{k+1} = P_k e_n, \quad u_{k+1} = P_k e_1.$$

证明 因为  $A_{k+1} = P_k^* A P_k$ ,

故

$$P_k A_{k+1} = A P_k,$$

$$P_{k+1} L_{k+1} = P_k Q_{k+1} L_{k+1}$$

$$= P_k (A_{k+1} - \sigma_{k+1} I) = (A - \sigma_{k+1} I) P_k,$$

$$P_{k+1}L_{k+1}e_n = (A - \sigma_{k+1}I)P_k e_n,$$

$$l_{nn}^{(k+1)}P_{k+1}e_n = (A - \sigma_{k+1}I)P_k e_n,$$

因为  $\|P_{k+1}e_n\|=1$ , 为了与(4)相比较, 将上式写成

$$P_{k+1}e_n = (A - \sigma_{k+1}I)P_k e_n / l_{nn}^{(k+1)}, \quad (6)$$

又由

$$P_1e_n = Q_1e_n = (A - \sigma_1I)e_n / l_{nn}^{(1)} = (A_1 - \sigma_1I)v_1 / l_{nn}^{(1)},$$

因为  $\|P_1e_n\|=1$ , 故  $v_2 = P_1e_n$ , 再从(6)就知道一般

$$v_{k+1} = P_k e_n,$$

且知  $l_{nn}^{(k)} = \nu_k$ .

另外  $L_{k+1}^* P_{k+1}^* = P_k^* (A - \sigma_{k+1}I),$

或

$$P_k L_{k+1}^* = (A - \sigma_{k+1}I)P_{k+1},$$

$$P_k L_{k+1}^* e_1 = (A - \sigma_{k+1}I)P_{k+1}e_1,$$

$$l_{11}^{(k+1)}P_k e_1 = (A - \sigma_{k+1}I)P_{k+1}e_1,$$

易知

$$u_{k+1} = P_k e_1, \quad l_{11}^{(k)} = \tau_k. \quad \text{证毕.}$$

从这个定理可知当  $\sigma_k = \sigma = \text{const}$  时,  $P_k$  的第一列常收敛到使  $1/|\lambda_i - \sigma|$  达到最大的那个特征值  $\lambda_i$  所对应的特征向量, 而且从

$$A_{k+1} = P_k^* A P_k,$$

可知  $A_{k+1}$  的元素  $a_{11}^{(k+1)} \rightarrow \lambda_i$ ,  $a_{1l}^{(k+1)} = a_{1l}^{(k+1)} \rightarrow 0, l=2, 3, \dots, n$ .

同样可知  $P_k$  的第  $n$  列, 当  $\sigma_k = \sigma = \text{const}$  时, 收敛到使  $|\lambda_i - \sigma|$  最大的特征值  $\lambda_i$  所对应的那个特征向量, 也有  $a_{nn}^{(k+1)} \rightarrow \lambda_i$ ,  $a_{ln}^{(k+1)} = a_{ln}^{(k+1)} \rightarrow 0, l=1, 2, \dots, n-1$ .

上述推论都是假定使  $1/|\lambda_i - \sigma|$  达到最大的特征值和使  $|\lambda_i - \sigma|$  达到最大的特征值都只有一个, 所对应的线性无关的特征向量也只有一个, 才获得的. 从乘幂法和反迭代法得到的结果, 收敛速度都是线性的, 因此并没有显示 QL 方法的内在优点.

## § 2 用于对称三对角矩阵时的 QL 方法的性质

我们的兴趣是在将 QL 方法用在对称三对角矩阵上, 因此我们要了解 QL 方法对于对称三对角矩阵有那些重要性质. 下面记

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \\ 0 & & \ddots & \beta_{n-1} & \\ & & & \beta_{n-1} & \alpha_n \end{pmatrix}$$

是对称三对角矩阵.

$$T - \sigma I = QL,$$

$$\hat{T} = LQ + \sigma I = \begin{pmatrix} \hat{\alpha}_1 & \hat{\beta}_1 & & & 0 \\ \hat{\beta}_1 & \hat{\alpha}_2 & \hat{\beta}_2 & & \\ & \hat{\beta}_2 & \ddots & \ddots & \\ 0 & & \ddots & \hat{\beta}_{n-1} & \\ & & & \hat{\beta}_{n-1} & \hat{\alpha}_n \end{pmatrix}.$$

**定理 3.5** 若  $T$  不可约, 则  $L$  的对角元  $l_{22}, l_{33}, \dots, l_{nn}$  全不为 0.

**证明** 因为  $T$  不可约, 故  $\beta_1, \beta_2, \dots, \beta_{n-1}$  全不为 0.  $T - \sigma I$  的第  $n$  列

$$t_n = (0, 0, \dots, 0, \beta_{n-1}, \alpha_n - \sigma)^T \neq 0,$$

故  $l_{nn} = \|t_n\| \neq 0,$

$T$  的第  $n-1$  列  $t_{n-1} = (0, 0, \dots, 0, \beta_{n-2}, \alpha_{n-1} - \sigma, \beta_{n-1})^T$ , 而

$$t_{n-1} - l_{n,n-1}q_n$$

的第  $n-2$  个元素为  $\beta_{n-2} \neq 0$ , 故

$$l_{n-1,n-1} = \|t_{n-1} - l_{n,n-1}q_n\| \neq 0,$$

以此类推可得  $l_{n-2,n-2} \neq 0$ ,  $l_{n-3,n-3} \neq 0$ ,  $\dots$ ,  $l_{33} \neq 0$ ,  $l_{22} \neq 0$ . 证毕.

**推论**  $T$  不可约,  $l_{11} = 0$  的充要条件是:  $\sigma$  是  $T$  的特征值.

定理 3.5 和它的推论加强了定理 3.1 的结果.

**定理 3.6** 若  $T$  不可约, 则

$$\hat{\beta}_i = \frac{l_{ii}}{l_{i+1,i+1}} \beta_i. \quad (7)$$

**证明**

$$q_n = \frac{1}{l_{nn}} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \beta_{n-1} \\ \alpha_n - \sigma \end{pmatrix},$$

因此  $q_n$  的第  $n-1$  个分量为  $\beta_{n-1}/l_{nn}$ , 第 1 个到第  $n-2$  个分量全为 0. 由

$$l_{kk}q_k = t_k - \sum_{j=k+1}^n l_{kj}q_j, \quad k=2, \dots, n-1,$$

可知  $q_k$  的第  $k-1$  个分量为  $\beta_{k-1}/l_{kk}$ , 第 1 个到第  $k-2$  个分量全为 0.

由  $\hat{T} = LQ + \sigma I$ ,

$\hat{T}$  的第  $i$  行第  $i+1$  列为  $\hat{\beta}_i$  是  $L$  的第  $i$  行与  $Q$  的第  $i+1$  列  $q_{i+1}$  相乘, 乘积为  $l_{ii}\beta_i/l_{i+1,i+1}$ , 故

$$\hat{\beta}_i = \frac{l_{ii}}{l_{i+1,i+1}} \beta_i, \quad i=1, 2, \dots, n-1. \text{ 证毕.}$$

推论 1  $T$  不可约, 且  $\sigma$  不是  $T$  的特征值, 则  $\hat{T}$  也是不可约.

推论 2  $T$  不可约, 则  $\hat{\beta}_2, \hat{\beta}_3, \dots, \hat{\beta}_{n-1}$  不为 0, 而  $\hat{\beta}_1=0$  的充要条件是  $\sigma$  为  $T$  的特征值.

从(7)知  $\hat{\beta}_1$  是  $\sigma$  的连续函数, 因此当  $\sigma$  充分接近特征值时,  $|\hat{\beta}_1|$  就会很小. 而当  $|\hat{\beta}_1|$  很小时,  $\hat{\alpha}_1$  就接近  $T$  的一个特征值. 后面这句话可以证明如下:

记  $T$  的特征值为  $\theta_1, \theta_2, \dots, \theta_n$ , 由相似性  $\hat{T}$  的特征值也是  $\theta_1, \theta_2, \dots, \theta_n$ . 记

$$\tilde{T} = \begin{pmatrix} \hat{\alpha}_1 & 0 & & & \\ 0 & \hat{\alpha}_2 & \hat{\beta}_2 & & \\ & \hat{\beta}_2 & \ddots & \ddots & \\ & & & \hat{\beta}_{n-1} & \hat{\alpha}_{n-1} \\ & & & \hat{\beta}_{n-1} & \hat{\alpha}_n \end{pmatrix},$$

它有一个特征值  $\hat{\alpha}_1$ . 而

$$\|\hat{T} - \tilde{T}\| = \left\| \begin{pmatrix} 0 & \hat{\beta}_1 & & & 0 \\ \hat{\beta}_1 & 0 & 0 & & \\ & 0 & \ddots & \ddots & 0 \\ & & & 0 & 0 \\ 0 & & & 0 & 0 \end{pmatrix} \right\| = |\hat{\beta}_1|,$$

由对称矩阵的特征值摄动定理[文献 9, p. 462], 知可找到  $\hat{T}$  的一个特征值  $\theta_i$ , 使

$$|\hat{\alpha}_1 - \theta_i| \leq \|\hat{T} - \tilde{T}\| = |\hat{\beta}_1|. \quad (8)$$

利用 Wielandt-Hoffman 定理[文献 9, p. 452]更有

$$(\hat{\alpha}_1 - \theta_i)^2 + (\hat{\beta}_1 - \theta_n)^2 + \dots + (\hat{\beta}_{n-1} - \theta_{n-1})^2 \leq 2\beta_1, \quad (9)$$

这里  $\tilde{\theta}_1, \dots, \tilde{\theta}_{n-1}$  是  $\hat{T}_{2,n}$  的特征值,  $i_1, i_2, \dots, i_{n-1}$  是  $1, 2, \dots, i-1, i+1, \dots, n$  的一种排列.

(8)、(9)两式告诉我们, 对于对称三对角矩阵  $T$ , 使用带位移  $\{\sigma_k\}$  的 QL 方法, 求特征值, 得到  $\{T^{(k)}\}$ , 如果  $T^{(k)}$  的  $\beta_1$  充分小时, 它的  $\alpha_1$  即为一个近似特征值, 接着只要对  $T_{2,n}^{(k)}$  继续进行 QL 方法, 求其余的特征值. 这样每求一个特征值, 还可把矩阵降低一阶, 求到的特征值还可以用 (8) 式和 (9) 式来估计误差, 这就显示出 QL 方法的优点. 这里也可看出 (8) 和 (9) 那样的估计式的意义. 如果能改进 (8) 式、(9) 式的结果, 是很有实用价值的.

W. Kahan 在 1966 年曾考虑过如下的问题: 对称三对角矩阵  $T$  中某一个  $\beta_i$  用 0 代替后的矩阵为  $\tilde{T}$ , 它的特征值为  $\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_n$ , 问  $\sum_{i=1}^n (\theta_i - \tilde{\theta}_i)^2$  能否得到比 (9) 更好一点的估计式. 他给出了一个结果

$$\sum_{i=1}^n (\theta_i - \tilde{\theta}_i)^2 \leq \frac{\beta_i^2}{\delta_i^2 + \rho_i^2} \left[ 2\rho_i^2 + \frac{\delta_i^2 \beta_i^2}{\delta_i^2 + \rho_i^2} \right],$$

其中

$$\delta_i = (\alpha_{i+1} - \alpha_i) / 2, \quad \rho_i^2 = (1 - 1/\sqrt{2}) (\beta_{i-1}^2 + \beta_{i+1}^2),$$

见文献 [10, p. 134].

对于 (8) 式至少可有

**定理 3.7** 对称三对角矩阵  $T$  的特征值为  $\theta_1, \theta_2, \dots, \theta_n$ ,

$$\theta_1 \leq \theta_2 \leq \theta_3 \leq \dots \leq \theta_n,$$

$T$  的某个  $\beta_i$  用 0 代替后的新矩阵  $\tilde{T}$ , 它的特征值为  $\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_n$ ,  $\tilde{\theta}_1 \leq \tilde{\theta}_2 \leq \dots \leq \tilde{\theta}_n$ , 假定  $\beta_{i-1}^2 + \beta_{i+1}^2 \neq 0$ , 则

$$|\theta_i - \tilde{\theta}_i| < |\beta_i|. \quad (10)$$

**证明** 不失一般性假定  $\beta_i > 0$ , 令

$$T = \tilde{T} + E,$$

$$E = \begin{pmatrix} 0 & & & 0 \\ & \ddots & & \\ 0 & 0 & \beta_i & \\ & \beta_i & \ddots & 0 \\ & & 0 & \ddots & 0 \\ & & & \ddots & 0 & 0 \end{pmatrix} \quad \begin{array}{l} \text{第 } i \text{ 行} \\ \text{第 } i+1 \text{ 行} \end{array}$$

易知  $E$  的特征值依小到大次序为  $-\beta_i, 0, 0, \dots, 0, \beta_i$ , 记对应的特征向量为  $b_1, b_2, \dots, b_n$ , 易知

$$b_1 = \frac{1}{\sqrt{2}}(e_i - e_{i+1}),$$

$$b_n = \frac{1}{\sqrt{2}}(e_i + e_{i+1}),$$

记  $\tilde{T}$  的特征值  $\tilde{\theta}_j$  对应的特征向量为  $x_j$ . 利用极大极小原理 (见 [9] 或本书第 4 章定理 4.1) 有

$$\theta_j = \max_{V_{n-j+1}} \min_{x \in V_{n-j+1}} \rho(x, T),$$

这里  $\rho(x, B) = (Bx, x) / (x, x)$ ,  $V_{n-j+1}$  是一个  $n-j+1$  维的线性子空间. 于是

$$\theta_j \geq \min_{x \in \{x_j, x_{j+1}, \dots, x_n\}} \rho(x, T),$$

右边的极小是能达到的, 记在  $\tilde{x}$  处达到, 可设  $\|\tilde{x}\| = 1$ ,

$$\theta_j \geq \rho(\tilde{x}, T) = \rho(\tilde{x}, \tilde{T}) + \rho(\tilde{x}, E).$$

对于  $\tilde{x}$  可分如下两种情况:

1. 若  $\tilde{x}$  是  $\tilde{T}$  的对应特征值  $\tilde{\theta}_j$  的特征向量, 此时

$$\rho(\tilde{x}, \tilde{T}) = \tilde{\theta}_j,$$

下面证明此时有  $\rho(\tilde{x}, E) > -\beta_i$ , 实际上因为  $\tilde{x} \neq \pm b_1$ , 否则  $b_1$  是  $\tilde{T}$  的对应  $\tilde{\theta}_j$  的特征向量, 但



$$(\tilde{T} - \tilde{\theta}_j) \mathbf{b}_1$$

$$= \frac{1}{\sqrt{2}} (0, \dots, 0, \underset{\text{第 } i-1 \text{ 列}}{\beta_{i-1}}, *, *, \underset{\text{第 } i+2 \text{ 列}}{-\beta_{i+1}}, 0, 0, \dots, 0)^T \neq 0,$$

因此  $(\tilde{\mathbf{x}}, \mathbf{b}_1)^2 = h_1^2 < 1$ .

$$\text{若 } \tilde{\mathbf{x}} = \sum_{i=1}^n h_i \mathbf{b}_i,$$

$$\rho(\tilde{\mathbf{x}}, E) = -\beta_i h_1^2 + h_n^2 \beta_i > -\beta_i.$$

2.  $\tilde{\mathbf{x}}$  不是  $\tilde{T}$  的对应  $\tilde{\theta}_j$  的特征向量, 此时

$$\rho(\tilde{\mathbf{x}}, \tilde{T}) > \tilde{\theta}_j, \quad \rho(\tilde{\mathbf{x}}, E) \geq -\beta_i,$$

因此不管哪种情况都有

$$\theta_j > \tilde{\theta}_j - \beta_i.$$

又利用

$$\theta_j = \min_{V_j} \max_{\mathbf{x} \in V_j} \rho(\mathbf{x}, T) \leq \max_{\mathbf{x} \in \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_j\}} \rho(\mathbf{x}, T),$$

同样可证得:  $\theta_j < \tilde{\theta}_j + \beta_i$ , 故

$$|\theta_j - \tilde{\theta}_j| < \beta_i.$$

从估计式(8)、(9)、(10)可知  $\{T^{(k)}\}$  中  $T^{(k)}$  的  $\beta_1$  趋于 0 的速度快, 那么  $\alpha_1$  趋于  $T$  的某个特征值的速度也快.

下面我们再给出一个估计  $\hat{\beta}_1$  的关系式(见文献[10]).

**定理 3.8** 设  $T - \sigma I = QL$ ,  $\hat{T} = LQ + \sigma I$ ,  $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$  是正交阵,  $L$  是下三角阵, 记  $\mathbf{q}_1$  与  $\mathbf{e}_1$  的夹角为  $\theta$ , 则

$$(T - \sigma I) \mathbf{q}_1 = l_{11} \mathbf{e}_1, \quad (11)$$

$$|\hat{\beta}_1| = l_{11} |\sin \theta|. \quad (12)$$

证明  $T - \sigma I = QL$ ,

$$T - \sigma I = L^T Q^T,$$

$$(T - \sigma I) Q = L^T,$$

故

$$(T - \sigma I) \mathbf{q}_1 = l_{11} \mathbf{e}_1,$$

(11)式获证. 又

$$Q\hat{T}=TQ,$$

故

$$\hat{\alpha}_1 \mathbf{q}_1 + \hat{\beta}_1 \mathbf{q}_2 = T \mathbf{q}_1,$$

$$\begin{aligned}\hat{\beta}_1 \mathbf{q}_2 &= T \mathbf{q}_1 - \hat{\alpha}_1 \mathbf{q}_1 = T \mathbf{q}_1 - \mathbf{q}_1^T T \mathbf{q}_1 \mathbf{q}_1 = (I - \mathbf{q}_1 \mathbf{q}_1^T) T \mathbf{q}_1 \\ &= (I - \mathbf{q}_1 \mathbf{q}_1^T) (l_{11} \mathbf{e}_1 + \sigma \mathbf{q}_1) \\ &= l_{11} (I - \mathbf{q}_1 \mathbf{q}_1^T) \mathbf{e}_1 = l_{11} (\mathbf{e}_1 - \mathbf{q}_1 \cos \theta),\end{aligned}$$

$$\text{于是} \quad \hat{\beta}_1^2 = l_{11}^2 (1 - 2 \cos^2 \theta + \cos^2 \theta) = l_{11}^2 \sin^2 \theta,$$

故(12)成立. 证毕.

推论  $|\hat{\beta}_1| \leq l_{11}$ .

由定理 3.8 可知, 为了估计  $\beta_1$  的收敛速度, 估计  $l_{11}$  是很重要的一环. 我们利用第 2 章定理 2.6 的结果给出如下的  $l_{11}$  和  $\hat{\beta}_1$  的表示式, 参见[31].

定理 3.9 设  $T$  不可约, 则

$$l_{11}^2 = \det(T - \sigma I)^2 / \omega_1, \quad (13)$$

其中

$$\begin{aligned}\omega_1 &= \det(T_{2,n} - \sigma I)^2 + (\det(T_{3,n} - \sigma I) \beta_1)^2 \\ &\quad + (\det(T_{4,n} - \sigma I) \beta_1 \beta_2)^2 + \cdots + (\beta_1 \beta_2 \cdots \beta_{n-1})^2.\end{aligned} \quad (14)$$

证明 先设  $\sigma$  不是  $T$  的特征值, 故  $\det(T - \sigma I) \neq 0$ , 由(11)  $T - \sigma I$  是一个对称三对角阵, 因此  $\mathbf{q}_1$  的分量  $q_{l1}$  ( $l=1, 2, \dots, n$ ), 可以利用定理 2.6 有表示式

$$q_{l1} = (-1)^{l-1} l_{11} \frac{\beta_1 \beta_2 \cdots \beta_{l-1} \det(T_{l+1,n} - \sigma I)}{\det(T - \sigma I)}, \quad (15)$$

但  $\sum_{i=1}^n q_{i1}^2 = 1$ , 故

$$\begin{aligned}l_{11}^2 (\det(T_{2,n} - \sigma I)^2 + (\det(T_{3,n} - \sigma I) \beta_1)^2 + \cdots \\ + (\beta_1 \beta_2 \cdots \beta_{n-1})^2) / \det(T - \sigma I)^2 = 1,\end{aligned}$$

因为  $T$  不可约, 故  $\omega_1 > 0$ , 于是即得(13).

如果  $\sigma$  是  $T$  的特征值, 则知  $l_{11}=0$ ,  $\det(T-\sigma I)=0$ , (13) 自然成立. 证毕.

**定理 3.10** 设  $T$  不可约, 令  $\omega_1$  如定理 3.9 中所定义,

$$\omega_2 = \det(T_{3,n} - \sigma I)^2 + (\det(T_{4,n} - \sigma I) \beta_2)^2 + \dots + (\beta_2 \beta_3 \dots \beta_{n-1})^2,$$

则

$$\hat{\beta}_1^2 = \beta_1^2 \det(T - \sigma I)^2 \omega_2 / \omega_1^2. \quad (16)$$

**证明** 由定理 3.8 知

$$\hat{\beta}_1^2 = l_{11}^2 \sin^2 \theta,$$

但  $q_{11} = \cos \theta$ ,  $\sin^2 \theta = 1 - q_{11}^2 = \sum_{i=2}^n q_{1i}^2$ , 故由 (15)

$$\begin{aligned} \sin^2 \theta &= l_{11}^2 \sum_{i=2}^n (\beta_1 \beta_2 \dots \beta_{i-1} \det(T_{i+1,n} - \sigma I))^2 / \det(T - \sigma I)^2 \\ &= l_{11}^2 \beta_1^2 \omega_2 / \det(T - \sigma I)^2 = \beta_1^2 \omega_2 / \omega_1, \end{aligned}$$

故 
$$\begin{aligned} \hat{\beta}_1^2 &= l_{11}^2 \sin^2 \theta = \frac{\det(T - \sigma I)^2}{\omega_1} \times \frac{\beta_1^2 \omega_2}{\omega_1} \\ &= \beta_1^2 \det(T - \sigma I)^2 \omega_2 / \omega_1^2. \text{ 证毕.} \end{aligned}$$

(16) 式给 QL 方法的收敛速度估计提供了非常有效的公式. 由此可知  $\hat{\beta}_1$  收敛于 0 的速度, 依赖于  $\det(T - \sigma I)$  收敛于 0 的速度. 对于不同的位移  $\sigma_k$ , 就有不同的  $\det(T - \sigma_k I)$  收敛状况, 本章的最后两节要介绍两种最常用的位移  $\{\sigma_k\}$  的取法, 并分析带那种位移时, QL 方法的收敛状况.

在研究收敛性前, 我们先介绍一下: 用旋转阵逐次相乘来实现 QL 迭代的计算方法; 它可以有递推的计算公式, 并且可以节省计算量. 对于对称三对角矩阵  $T$ , 和位移  $\sigma$ ,  $T - \sigma I = QL$ ,  $\hat{T} = LQ + \sigma I$ , 我们要给出从  $T$  求  $\hat{T}$  的计算公式, 也即从  $\alpha_1, \alpha_2, \dots, \alpha_n, \beta_1, \beta_2, \dots, \beta_{n-1}$  求  $\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_n, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_{n-1}$  的计算公式. 记  $\alpha_k = \alpha_k - \sigma$ ,  $k=1, 2, \dots, n$ , 即

$$T - \sigma I = \begin{pmatrix} \bar{\alpha}_1 & \beta_1 & & & 0 \\ \beta_1 & \bar{\alpha}_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \\ 0 & & \ddots & \beta_{n-1} & \alpha_n \\ & 0 & & \beta_{n-1} & \alpha_n \end{pmatrix},$$

将它左乘旋转阵  $R_{n-1, n}$ , 将  $(n-1, n)$  位置上的元素化成零.

$$\begin{pmatrix} 1 & & & & 0 \\ & \ddots & & & \\ & & 1 & & \\ & & c_n & -s_n & \\ 0 & s_n & c_n & & \end{pmatrix} \begin{pmatrix} \bar{\alpha}_1 & \beta_1 & & & 0 \\ \beta_1 & \bar{\alpha}_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \beta_{n-2} \\ & & \ddots & \beta_{n-2} & \bar{\alpha}_{n-1} & \beta_{n-1} \\ 0 & & & \beta_{n-1} & \bar{\alpha}_n \end{pmatrix} \\ = \begin{pmatrix} \bar{\alpha}_1 & \beta_1 & & & \\ \beta_1 & \bar{\alpha}_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \beta_{n-3} \\ & & \ddots & \beta_{n-3} & \bar{\alpha}_{n-2} & \beta_{n-2} \\ & & & c_n \bar{\alpha}_{n-1} - s_n \beta_{n-1} & c_n \beta_{n-1} - s_n \bar{\alpha}_n \\ & & & s_n \bar{\alpha}_{n-1} + \beta_{n-1} c_n & s_n \beta_{n-1} + c_n \bar{\alpha}_n \end{pmatrix},$$

为了使  $(n-1, n)$  位置上的元素为 0, 只要  $c_n \beta_{n-1} - s_n \bar{\alpha}_n = 0$ , 即得

$$c_n = \bar{\alpha}_n / r_n, \quad s_n = \beta_{n-1} / r_n, \quad r_n = \sqrt{\bar{\alpha}_n^2 + \beta_{n-1}^2},$$

记  $t_n = s_n \beta_{n-2}$ ,  $q_n = s_n \bar{\alpha}_{n-1} + \beta_{n-1} c_n$ ,  $p_{n-1} = c_n \bar{\alpha}_{n-1} - s_n \beta_{n-1}$ ,



将(17)式与前面  $R_{n-1,n}$  旋转时的量进行对比, 就会发现, 如果令  $\bar{\alpha}_n = p_n$ ,  $c_{n+1} = 1$ , 就有一般的递推公式:

$$\begin{cases} r_k = (\beta_{k-1}^2 + p_k^2)^{\frac{1}{2}}, \\ s_k = \beta_{k-1}/r_k, \quad c_k = p_k/r_k, \\ p_{k-1} = c_k \bar{\alpha}_{k-1} - s_k c_{k+1} \beta_{k-1}, \\ q_k = s_k \bar{\alpha}_{k-1} + c_k c_{k+1} \beta_{k-1}, \\ t_k = s_k \beta_{k-2}, \quad k = n, n-1, \dots, 2, \end{cases} \quad (18)$$

得到

$$R_{1,2} R_{2,3} \cdots R_{n-2,n-1} R_{n-1,n} (T - \sigma I) = \begin{pmatrix} p_1 & & & & 0 \\ & q_2 & r_2 & & \\ & t_3 & q_3 & r_3 & \\ & & \ddots & \ddots & \ddots \\ & & & 0 & t_n & q_n & r_n \end{pmatrix} = L.$$

同时知道

$$Q = R_{n-1,n}^* R_{n-2,n-1}^* \cdots R_{2,3}^* R_{1,2}^*.$$

现在来考察  $LQ = LR_{n-1,n}^* R_{n-2,n-1}^* \cdots R_{2,3}^* R_{1,2}^*$ ,  $LR_{n-1,n}^*$  只改变  $L$  的最后两列; 而  $LR_{n-1,n}^*$  再右乘  $R_{n-2,n-1}^* \cdots R_{2,3}^* R_{1,2}^*$ , 不改变最后一列, 因此从  $LR_{n-1,n}^*$  的最后一列可以得到  $\hat{\beta}_{n-1}$  和  $\hat{\alpha}_n$ . 通过简单的相乘得到

$$\begin{cases} \hat{\beta}_{n-1} = r_{n-1} s_n, \\ \hat{\alpha}_n = q_n s_n + r_n c_n + \sigma, \end{cases} \quad (19)$$

$$LR_{n-1,n}^* = \begin{pmatrix} p_1 & & & & & & 0 \\ & q_2 & r_2 & & & & \\ & & \ddots & \ddots & & & \\ & t_3 & & \ddots & & & \\ & & & & r_{n-2} & & \\ & & & & & q_{n-1} & r_{n-1}c_n & \hat{\beta}_{n-1} \\ 0 & & & & & & & \\ & & & & t_n & x & q_n s_n + r_n c_n & \end{pmatrix},$$

对  $LR_{n-1,n}^*$  再右乘  $R_{n-2,n-1}^*$ , 只改变第  $n-2$  列和  $n-1$  列, 得到

$$\begin{cases} \hat{\beta}_{n-2} = r_{n-2}s_{n-1}, \\ \hat{\alpha}_{n-1} = q_{n-1}s_{n-1} + r_{n-1}c_n c_{n-1} + \sigma, \end{cases} \quad (20)$$

在  $LR_{n-1,n}^* R_{n-2,n-1}^*$  中  $(n-2, n-2)$  位置上的元素变成  $c_{n-1}r_{n-2}$ , 比较 (19) 与 (20), 如果考虑到  $c_{n+1}=1$ , 再令  $p_1=r_1$ , 则有一般的递推公式

$$\left. \begin{aligned} \hat{\beta}_k &= r_k s_{k+1}, \\ \hat{\alpha}_{k+1} &= q_{k+1} s_{k+1} + r_{k+1} c_{k+2} c_{k+1} + \sigma, \\ k &= n-1, n-2, \dots, 1, \\ \hat{\alpha}_1 &= r_1 c_2 + \sigma, \end{aligned} \right\}$$

考虑到  $r_{k+1}c_{k+1}=p_{k+1}$ , 上式可以写成

$$\left\{ \begin{aligned} \hat{\beta}_k &= r_k s_{k+1}, \\ \hat{\alpha}_{k+1} &= q_{k+1} s_{k+1} + p_{k+1} c_{k+2} + \sigma, \\ k &= n-1, n-2, \dots, 1, \\ \hat{\alpha}_1 &= r_1 c_2 + \sigma. \end{aligned} \right. \quad (21)$$

(18)式与(21)式合并在一起, 就是从  $T$  计算  $\hat{T}$  的计算公式, 在(18)式中的  $t_k$  这个量, 在(21)式中没有用到, 在(18)式中的其他量中也没有用到它, 因此可以不必计算。这样总共需:

加减法运算  $5n-4$  次,  
 乘法运算  $10n-9$  次,  
 除法运算  $2n-2$  次,  
 开方运算  $n-1$  次.

一般计算机对开方运算总是多费一点时间, 因此很多数值分析专家要想法把开方运算去掉, 对此提出不少算法, 其中 C. H. Reinsch 在 1971 年发表的文章 [11] 中的算法是非常成功的. 他是对 QR 方法来说的, 我们将它修改成对 QL 方法使用的, 介绍如下: 在算式 (18) 和 (21) 中引进二个新的量

$$h_k = p_k c_{k+1}, \quad g_k = p_k / c_{k+1},$$

有  $h_n = g_n = p_n = \bar{\alpha}_n$ , 再注意到成立等式

$$c_k \beta_{k-1} = s_k p_k,$$

于是从 (18) 式

$$p_{k-1} = c_k \bar{\alpha}_{k-1} - s_k c_{k+1} \beta_{k-1},$$

可得

$$\begin{aligned}
 g_{k-1} &= \bar{\alpha}_{k-1} - s_k c_{k+1} \beta_{k-1} / c_k = \bar{\alpha}_{k-1} - s_k c_{k+1} \beta_{k-1}^2 / c_k \beta_{k-1} \\
 &= \bar{\alpha}_{k-1} - s_k c_{k+1} \beta_{k-1}^2 / s_k p_k = \bar{\alpha}_{k-1} - \beta_{k-1}^2 / g_k,
 \end{aligned} \quad (22)$$

这是一个从  $g_k$  算  $g_{k-1}$  的递推公式, 式中没有开方运算. 又

$$h_{k-1} = g_{k-1} p_k^2 / r_k^2, \quad (23)$$

也即知道  $p_k, r_k$  后可以算  $h_{k-1}$ . 而

$$r_k^2 = p_k^2 + \beta_{k-1}^2, \quad (24)$$

$$p_{k-1}^2 = g_{k-1} h_{k-1}, \quad (25)$$

这些都是没有开方运算的算式, 它们告诉我们从  $T$  可以不经开方运算, 得到  $r_k, p_k, h_k, g_k$ , 因此下面只要将  $\hat{\beta}_k$  和  $\hat{\alpha}_k$  用这些量和  $T$  的元素表示出来, 表示式中要求不包含开方运算. 从 (21)



$$\begin{aligned}
\hat{\alpha}_{k+1} &= q_{k+1}s_{k+1} + p_{k+1}c_{k+2} + \sigma \\
&= h_{k+1} + \sigma + (s_{k+1}\bar{\alpha}_k + c_{k+1}c_{k+2}\beta_k)s_{k+1} \\
&= h_{k+1} + \sigma + s_{k+1}^2(\bar{\alpha}_k + h_{k+1}), \\
\hat{\alpha}_1 &= r_1c_2 + \sigma = h_1 + \sigma, \\
\hat{\beta}_k^2 &= r_k^2s_{k+1}^2,
\end{aligned}$$

综合这些式子得到无开方运算的 QL 迭代的算式,

$$\text{初始} \quad p_n = g_n = h_n = \bar{\alpha}_n, \quad r_n^2 = p_n^2 + \beta_{n-1}^2,$$

$$\left\{ \begin{array}{l}
s_k^2 = \beta_{k-1}^2 / r_k^2, \\
\hat{\alpha}_k = h_k + \sigma + s_k^2(\bar{\alpha}_{k-1} + h_k), \\
g_{k-1} = \bar{\alpha}_{k-1} - \beta_{k-1}^2 / g_k, \\
h_{k-1} = g_{k-1}p_k^2 / r_k^2, \quad p_{k-1}^2 = g_{k-1}h_{k-1}, \\
r_{k-1}^2 = p_{k-1}^2 + \beta_{k-2}^2, \quad \beta_0^2 = 0, \\
\hat{\beta}_{k-1}^2 = r_{k-1}^2s_k^2, \\
k = n, n-1, \dots, 2, \\
\hat{\alpha}_1 = h_1 + \sigma.
\end{array} \right. \quad (26)$$

这是从  $\alpha_1, \alpha_2, \dots, \alpha_n, \beta_1^2, \beta_2^2, \dots, \beta_{n-1}^2$  算出  $\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_n, \hat{\beta}_1^2, \hat{\beta}_2^2, \dots, \hat{\beta}_{n-1}^2$  的公式. 尽管算出的不是  $\hat{\beta}_k$  而是  $\hat{\beta}_k^2$ , 但是对 QL 方法来说, 这是没有妨害的.

不计从  $\beta_1, \beta_2, \dots, \beta_{n-1}$  得到  $\beta_1^2, \beta_2^2, \dots, \beta_{n-1}^2$  的计算量, 总的计算量是

加减法运算:  $5n-3$  次,

除法运算:  $3n-3$  次,

乘法运算:  $4n-3$  次.

在计算过程中遇到  $g_i \simeq 0$  时, 可以用小的正数  $\delta$  来代替, 计算可以继续, 从  $g_i$  表示式中可知这样的代替相当于对  $\alpha_i$  有一个摄动, 变成  $\alpha_i + \delta$ ; 只要  $\delta$  充分小, 仍可使求得特征值满足预定的精度要求.

### § 3 带 Rayleigh 商位移的 QL 方法

对于带位移  $\{\sigma_k\}$  的 QL 方法:

$$T^{(k)} - \sigma_k I = Q_k L_k,$$

$$T^{(k+1)} = L_k Q_k + \sigma_k I,$$

如果

$$T^{(1)} = T,$$

$$T^{(k)} = \begin{pmatrix} \alpha_1^{(k)} & \beta_1^{(k)} & & 0 \\ \beta_1^{(k)} & \alpha_2^{(k)} & \beta_2^{(k)} & \\ & \beta_2^{(k)} & \ddots & \beta_{n-1}^{(k)} \\ 0 & & \beta_{n-1}^{(k)} & \alpha_n^{(k)} \end{pmatrix},$$

而取  $\alpha_1^{(k)}$  作为  $\sigma_k$ , 那么这时的 QL 方法, 称为带 Rayleigh 商位移的 QL 方法.

为了弄清为何叫做 Rayleigh 商位移, 先要了解什么是 Rayleigh 商. 对于任意对称矩阵  $A$ , 取向量  $\mathbf{x} \neq 0$ , 那么

$$\rho(\mathbf{x}) = (A\mathbf{x}, \mathbf{x}) / (\mathbf{x}, \mathbf{x})$$

称为  $A$  在向量  $\mathbf{x}$  时的 Rayleigh 商, 它当然跟  $A$  有关, 因此有时为了强调一下对应的矩阵, 记为  $\rho(\mathbf{x}, A)$ . 由对称矩阵的极大极小原理知道  $\rho(\mathbf{x})$  的极值就是  $A$  的特征值, 因此可以通过 Rayleigh 商来求得  $A$  的特征值.

所谓 Rayleigh 商迭代法, 就是一种求对称矩阵特征值和对应特征向量的方法. 算法如下:

1. 取一个单位向量  $\mathbf{x}_1$ ,  $1 \rightarrow k$ ;
2. 计算  $\rho_k = \rho(\mathbf{x}_k)$ ;
3. 如果  $\det(A - \rho_k I) \neq 0$ , 则从  $(A - \rho_k I)\mathbf{y}_{k+1} = \mathbf{x}_k$ , 求

$\mathbf{y}_{k+1}$ , 再将  $\mathbf{y}_{k+1}$  标准化成为  $\mathbf{x}_{k+1}$ ,  $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} / \|\mathbf{y}_{k+1}\|$ , 如果  $\|\mathbf{y}_{k+1}\|$  很大,  $\rho_k$  作为  $A$  的特征值近似值,  $\mathbf{x}_{k+1}$  作为对应特征向量的近似, 计算完成.

如果  $\|\mathbf{y}_{k+1}\|$  不很大, 则  $k+1 \rightarrow k$  转 2.

如果  $\det(A - \rho_k I) = 0$ , 从

$$(A - \rho_k I) \mathbf{x}_{k+1} = 0$$

算得单位向量  $\mathbf{x}_{k+1}$ ,  $\rho_k$ ,  $\mathbf{x}_{k+1}$  是  $A$  的特征值和对应的特征向量, 计算完成.

对于 Rayleigh 商迭代法, 是否从任意  $\mathbf{x}_1$  出发,  $\rho_k$  都是收敛到  $A$  的一个特征值? 这就是 Rayleigh 商迭代法的收敛性问题. 我们先不回答这个问题. 还是再来回答那样的位移, 为何叫带 Rayleigh 商位移. 有如下定理 [10].

**定理 3.11** 对  $T^{(1)} = T$  按带 Rayleigh 位移的 QL 方法获得矩阵序列  $\{T^{(k)}\}$ , 和从  $\mathbf{x}_1 = \mathbf{e}_1$  按 Rayleigh 迭代法获得的序列  $\{\rho_k\}$ , 成立如下关系

$$\rho_k = \alpha_1^{(k)}, \quad (27)$$

$$\mathbf{x}_{k+1} = P_k \mathbf{e}_1. \quad (28)$$

**证明** 因为

$$\rho_1 = (T \mathbf{x}_1, \mathbf{x}_1) / (\mathbf{x}_1, \mathbf{x}_1) = (T \mathbf{e}_1, \mathbf{e}_1) / (\mathbf{e}_1, \mathbf{e}_1) = \alpha_1^{(1)},$$

因此 (27) 式对  $k=1$  成立, 又

$$(T - \rho_1 I) \mathbf{y}_2 = \mathbf{x}_1 = \mathbf{e}_1,$$

由定理 3.8 知  $\mathbf{y}_2$  与  $\mathbf{q}_1^{(1)}$  至多差一个正常数因子, 即  $\mathbf{x}_2 = \mathbf{q}_1^{(1)} = Q_1 \mathbf{e}_1 = P_1 \mathbf{e}_1$ , 因此 (28) 式也对  $k=1$  成立.

现在用数学归纳法, 假定 (27)、(28) 对  $k=1, 2, \dots, j$  都成立, 考虑  $k=j+1$  的情况, 此时

$$\begin{aligned} \rho_{j+1} &= (T \mathbf{x}_{j+1}, \mathbf{x}_{j+1}) / (\mathbf{x}_{j+1}, \mathbf{x}_{j+1}) \\ &= (T P_j \mathbf{e}_1, P_j \mathbf{e}_1) = (P_j^* T P_j \mathbf{e}_1, \mathbf{e}_1), \end{aligned}$$

但由 QL 方法知

$$P_j^* T P_j = T^{(j+1)}$$

成立, 于是  $\rho_{j+1} = (T^{(j+1)} \mathbf{e}_1, \mathbf{e}_1) = \alpha_1^{(j+1)}$ ,

又  $(T - \rho_{j+1}) \mathbf{y}_{j+2} = \mathbf{x}_{j+1}$ ,

$$(T - \alpha_1^{(j+1)}) \mathbf{y}_{j+2} = P_j \mathbf{e}_1,$$

$$(P_j^* T P_j - \alpha_1^{(j+1)}) P_j^* \mathbf{y}_{j+2} = \mathbf{e}_1,$$

$$(T^{(j+1)} - \alpha_1^{(j+1)}) P_j^* \mathbf{y}_{j+2} = \mathbf{e}_1,$$

再由定理 3.8 知  $P_j^* \mathbf{y}_{j+1}$  与  $\mathbf{q}_1^{(j+1)}$  至多差一个正常数因子, 即

$$P_j^* \mathbf{y}_{j+2} = \tau \mathbf{q}_{j+1} \mathbf{e}_1,$$

或  $\mathbf{y}_{j+2} = \tau P_j \mathbf{q}_{j+1} \mathbf{e}_1 = \tau P_{j+1} \mathbf{e}_1$ ,

从而  $\mathbf{x}_{j+2} = P_{j+1} \mathbf{e}_1$ . 证毕.

根据这个定理可知, 带 Rayleigh 商位移的 QL 方法中的  $\beta_1$  是否收敛于零跟从  $\mathbf{e}_1$  出发的 Rayleigh 商迭代是否收敛, 是等价的. 实际上, 若  $\rho_k \rightarrow \theta_i$ ,  $\theta_i$  是  $T$  的一个特征值,  $\mathbf{x}_k \rightarrow \mathbf{y}_i$ ,  $\mathbf{y}_i$  是对应  $\theta_i$  的单位特征向量, 于是按定理 3.11,  $P_j$  的第 1 列趋于  $\mathbf{y}_i$ , 由

$$T^{(j+1)} = P_j^* T P_j = P_j^* [\theta_i \mathbf{y}_i + \varepsilon, T \mathbf{g}_2, \dots, T \mathbf{g}_n]$$

知  $\beta_1^{(j+1)} \rightarrow 0$ , 这里  $\varepsilon$  是趋于 0 的向量,  $\mathbf{g}_2, \dots, \mathbf{g}_n$  是正交矩阵  $P_j$  的第 2 列、 $\dots$ 、第  $n$  列. 反过来也成立, 如果  $\beta_1^{(k)} \rightarrow 0$ , 则  $\rho_k \rightarrow \theta_i$ ,  $\mathbf{x}_{k+1} \rightarrow \mathbf{y}_i$ .

自然人们会提出这样的问题: 对任意对称三对角矩阵  $T$ , 使用带 Rayleigh 商位移的 QL 方法,  $\beta_1^{(k)}$  是否一定能收敛于零? 下面的例子说明对有些矩阵  $T$ ,  $\beta_1^{(k)}$  是不收敛的.

$$\text{例 } T = \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} \text{ 使用带 Rayleigh 商位移的 QL 方法.}$$

$$T^{(1)} = T, \sigma_1 = \alpha_1^{(1)} = \alpha_1 = 0, T^{(1)} - \sigma_1 I = Q_1 L_1,$$

$$Q_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, L_1 = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix},$$

$$T^{(2)} = L_1 Q_1 + \sigma_1 I = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix},$$

由此可知  $T^{(k)} = T^{(1)}, k=2, 3, \dots, \sigma_k = \alpha_1^{(k)} = 0,$

$$\beta_1^{(k)} = \frac{1}{2}$$

不收敛于零。

这个例子当然是很特殊的情况。在这个情况下, Rayleigh 商迭代法, 取  $\mathbf{x}_1 = \mathbf{e}_1$ ,  $\rho_k$  也是不能收敛到特征值的。我们来考察一下, 此时  $\rho_k$  和  $\mathbf{x}_k$  的状况。

$$\mathbf{x}_1 = \mathbf{e}_1, \rho_1 = \rho(\mathbf{x}_1) = (T\mathbf{e}_1, \mathbf{e}_1) = 0, (T - \rho_1 I)\mathbf{y}_2 = \mathbf{x}_1,$$

即 
$$\begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} \begin{pmatrix} y_{12} \\ y_{22} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix},$$

得到 
$$\mathbf{y}_2 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix},$$

$$\rho_2 = \rho(\mathbf{x}_2) = (T\mathbf{x}_2, \mathbf{x}_2) = 0,$$

再由 
$$(T - \rho_2)\mathbf{y}_3 = \mathbf{x}_2,$$

即 
$$\begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} \begin{pmatrix} y_{13} \\ y_{23} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix},$$

得  $y_3 = (2, 0)^T$ ,  $x_3 = (1, 0)^T$ .

由此可知:  $\rho_k = 0$ ,  $x_{2k-1} = (1, 0)^T$ ,  $x_{2k} = (0, 1)^T$ . 尽管  $\rho_k$  始终是 0, 但是 0 不是特征值,  $T$  的特征值是  $\theta_1 = \frac{1}{2}$ ,  $\theta_2 = -\frac{1}{2}$ , 对应的特征向量  $y_1 = (1, 1)^T / \sqrt{2}$ ,  $y_2 = (1, -1)^T / \sqrt{2}$ . 恰有  $\rho_k = \theta_1 + \theta_2$ ,  $x_{2k-1} = (y_1 + y_2) / \sqrt{2}$ ,  $x_{2k} = (y_1 - y_2) / \sqrt{2}$ . 这一现象反映了 B. N. Parlett 和 W. Kahan 在 1969 年给出的一个定理.

定理: 任意取一个单位向量  $x_1$ , 作为初始向量, Rayleigh 商迭代法所产生的序列  $\rho_k$ ,  $x_k$ , 下列两者必有一个发生:

1.  $\rho_k$  收敛到一个特征值,  $x_k$  收敛到对应的特征向量. 收敛速度是三次的.

2.  $\rho_k$  收敛到一个数  $\bar{\rho}$ ,  $\bar{\rho}$  不是  $T$  的特征值,  $x_{2k}$  收敛到向量  $x_+$ ,  $x_{2k-1}$  收敛到向量  $x_-$ ,  $x_+ + x_-$  和  $x_+ - x_-$  是  $T$  的两个特征向量, 此时收敛速度是线性的.

有关这个定理的证明, 可参看 [10, 第四章], 不在此介绍了.

上述 Parlett 和 Kahan 定理指出 Rayleigh 商迭代收敛到特征值时, 收敛速度是三次的; 对于带 Rayleigh 商位移的 QL 方法, 如果  $\beta_1^{(k)} \rightarrow 0$ , 那么它的收敛速度也是三次的, 有如下定理:

**定理 3.12** 设  $T$  不可约, 它的特征值为  $\theta_1, \theta_2, \dots, \theta_n$  由带 Rayleigh 商位移的 QL 方法, 产生  $\{T^{(k)}\}$ , 如果  $\beta_1^{(k)} \rightarrow 0$ , 则  $\alpha_1^{(k)} \rightarrow \theta_i$ , 并且

$$|\beta_1^{(k+1)}| = |\beta_1^{(k)}|^3 |\det(T_{3,n}^{(k)} - \sigma_k I)| (\omega_2^{(k)})^{1/2} / |\omega_1^{(k)}|, \quad (29)$$

这里

$$\omega_1^{(k)} = \det(T_{2,n}^{(k)} - \sigma_k I)^2 + (\det(T_{3,n}^{(k)} - \sigma_k I) \beta_1)^2 + \cdots + (\beta_1^{(k)} \beta_2^{(k)} \cdots \beta_{n-1}^{(k)})^2,$$

有

$$\lim_{k \rightarrow \infty} \omega_1^{(k)} = \left( \prod_{\substack{l=1 \\ l \neq i}}^n (\theta_l - \theta_i) \right)^2 \neq 0, \quad (30)$$

$$\omega_2^{(k)} = \det(T_{3,n}^{(k)} - \sigma_k I)^2 + (\det(T_{4,n}^{(k)} - \sigma_k I) \beta_2)^2 + \cdots + (\beta_2^{(k)} \beta_3^{(k)} \cdots \beta_{n-1}^{(k)})^2$$

是一致有界量。

证明  $\beta_1^{(k)} \rightarrow 0$ , 因此  $\alpha_1^{(k)} = \sigma_k \rightarrow \theta_i$  (参见[31]),  $\theta_i$  是  $T$  的某个特征值。由此可知  $|\det(T_{i,n}^{(k)} - \sigma_k I)|$  是一致有界的, 另外  $|\beta_i^{(k)}| = |(T^{(k)} e_{i+1}, e_i)| \leq \|T^{(k)}\| = \|T\|$ ,  $i=1, 2, \dots, n-1$ , 因此  $\omega_2^{(k)}$  是一致有界量。

若  $T_{2,n}^{(k)}$  的特征值为  $\theta_1^{(k)}, \theta_2^{(k)}, \dots, \theta_{n-1}^{(k)}$ , 由估计式 (9) 知  $\theta_1^{(k)}, \theta_2^{(k)}, \dots, \theta_{n-1}^{(k)}$  收敛到  $\theta_i, \theta_{i+1}, \dots, \theta_{n-1}$ ,

$$\det(T_{2,n}^{(k)} - \sigma_k I) = \prod_{l=1}^{n-1} (\theta_l^{(k)} - \sigma_k),$$

故 
$$\lim_{k \rightarrow \infty} \det(T_{2,n}^{(k)} - \sigma_k I) = \prod_{\substack{l=1 \\ l \neq i}}^n (\theta_l - \theta_i),$$

由此可得 
$$\lim_{k \rightarrow \infty} \omega_1^{(k)} = \left( \prod_{\substack{l=1 \\ l \neq i}}^n (\theta_l - \theta_i) \right)^2,$$

因  $T$  不可约, 于是  $T$  的特征值全不相同, 故上式不等于 0。

利用定理 3.10, 有

$$|\beta_1^{(k+1)}| = |\beta_1^{(k)}| |\det(T^{(k)} - \sigma_k I)| (\omega_2^{(k)})^{\frac{1}{2}} / |\omega_1^{(k)}|, \quad (31)$$

但

$$\begin{aligned} \det(T^{(k)} - \sigma_k I) &= (\alpha_1^{(k)} - \sigma_k) \det(T_{2,n}^{(k)} - \sigma_k I) \\ &\quad - (\beta_1^{(k)})^2 \det(T_{3,n}^{(k)} - \sigma_k I), \end{aligned}$$

因为  $\alpha_1^{(k)} = \sigma_k$ , 故

$$\det(T^{(k)} - \sigma_k I) = -(\beta_1^{(k)})^2 \det(T_{3,n}^{(k)} - \sigma_k I),$$

将此代入(31)即得(29)式。证毕。

由定理 3.12 可知  $\beta_1^{(k)} \rightarrow 0$  时, 它的收敛速度至少是三次的。但是满足什么样条件的对称三对角矩阵, 可以保证上述  $\beta_1^{(k)} \rightarrow 0$ , 仍然是一个值得研究的问题。

#### § 4 带 Wilkinson 位移的 QL 方法

鉴于带 Rayleigh 商位移的 QL 方法, 偶而有  $\beta_1^{(k)}$  不收敛于 0 的情况发生, 目前用的较多的是带 Wilkinson 位移的 QL 方法。这种位移取法在 Wilkinson 1965 出版的书[20]中已经提到。1968 年 Wilkinson 在文[21]中证明: 这种取法的位移, 能保证对任意对称三对角矩阵,  $\beta_1^{(k)}$  都收敛于 0, 并证明收敛速度至少是二次的。但是证明比较复杂, 也可参看 C. L. Lawson 和 R. J. Hanson 的书[22, Appendix B]。

1978 年 Hoffman 和 Parlett 给出一个很巧妙的关于  $\beta_1^{(k)}$  收敛于 0 的证明, 同时也给出在一般情况下  $\hat{\beta}_1 \sim 0(\beta_1^3 \beta_2^2)$ , 在特殊情况下  $\hat{\beta}_1 \sim 0(\beta_1^2)$  的收敛速度估计。不过在证明中, 用到这样的命题: 对于方程

$$(T - \sigma I)p = e_1$$

的解  $p = (\pi_1, \pi_2, \pi_3, \dots, \pi_n)^T$ , 当  $\beta_1$  充分小时它的分量中模最大的是  $\pi_1$  和  $\pi_2$ 。这个命题没有严格的证明。B. N. Parlett 在 1980 年出版的书[10, p. 155~156]中, 在附加条件

$$\beta_2 \rightarrow 0, \quad \beta_3 \rightarrow 0, \quad \alpha_i \rightarrow \theta_i \quad (i=1, 2, 3), \quad (32)$$

证明了  $|\hat{\beta}_1 / \beta_1^3 \beta_2| \rightarrow 1 / |\theta_2 - \theta_1|^3 |\theta_3 - \theta_1| \neq 0$ 。

这里我们证明  $\beta_1^{(k)} \rightarrow 0$ , 是采用书[10]中的证明, 关于收敛速度的估计我们去掉了附加的假定(32), 给出  $\hat{\beta}_1$  与  $\beta_1$ 、



$\beta_2, \beta_3, \dots$  的正确关系式. 指出如果  $\lim_{k \rightarrow \infty} (\alpha_2^{(k)} - \sigma_k) \neq 0$ , 则  $\hat{\beta}_1 = o(\beta_1^3 \beta_2^2)$ . 任何情况下  $\hat{\beta}_1 = o(|\alpha_1 - \sigma| \beta_1 \beta_2^2)$ .

先介绍什么叫 Wilkinson 位移. 对于

$$T = \begin{pmatrix} \alpha_1 & \beta_1 & & 0 \\ & \alpha_2 & \beta_2 & \\ & \ddots & \ddots & \ddots \\ 0 & & \beta_{n-1} & \alpha_n \end{pmatrix},$$

它的二阶顺序主子阵  $\begin{pmatrix} \alpha_1 & \beta_1 \\ \beta_1 & \alpha_2 \end{pmatrix}$ , 有两个实根

$$\omega_1 = [\alpha_1 + \alpha_2 + \sqrt{(\alpha_1 - \alpha_2)^2 + 4\beta_1^2}] / 2,$$

$$\omega_2 = [\alpha_1 + \alpha_2 - \sqrt{(\alpha_1 - \alpha_2)^2 + 4\beta_1^2}] / 2,$$

取其中靠近  $\alpha_1$  的那个根, 作为位移值  $\sigma$ , 带这样取法的位移, 称为带 Wilkinson 位移. 当然对于  $T^{(k)}$  来说上述  $\alpha_1, \alpha_2, \beta_1$  和  $\sigma$  都要标以  $k$  的.

如果记  $\delta = (\alpha_2 - \alpha_1) / 2$ , 那么

$$\omega_1 = \frac{\alpha_1 + \alpha_2}{2} + \sqrt{\delta^2 + \beta_1^2} = \frac{\alpha_1 + \alpha_2}{2} - \delta + \delta + \sqrt{\delta^2 + \beta_1^2}$$

$$= \alpha_1 + \delta + \sqrt{\delta^2 + \beta_1^2} = \alpha_1 - \beta_1^2 / (\delta - \sqrt{\delta^2 + \beta_1^2}),$$

即  $\omega_1 - \alpha_1 = -\beta_1^2 / (\delta - \sqrt{\delta^2 + \beta_1^2}),$

同样  $\omega_2 - \alpha_1 = -\beta_1^2 / (\delta + \sqrt{\delta^2 + \beta_1^2}),$

因此当  $\delta > 0$  时

$$|\omega_2 - \alpha_1| < |\omega_1 - \alpha_1|,$$

当  $\delta < 0$  时

$$|\omega_1 - \alpha_1| < |\omega_2 - \alpha_1|.$$

因此当  $\delta \geq 0$  取  $\sigma = \omega_2$ , 当  $\delta < 0$  时取  $\sigma = \omega_1$ , 利用函数

$$\text{sign } \delta = \begin{cases} 1, & \delta \geq 0; \\ -1, & \delta < 0, \end{cases}$$

不管那种情况可取

$$\sigma = \alpha_1 - \text{sign } \delta (\beta_1^2 / (|\delta| + \sqrt{\delta^2 + \beta_1^2})). \quad (33)$$

因为  $\sigma$  是  $\begin{pmatrix} \alpha_1 & \beta_1 \\ \beta_1 & \alpha_1 \end{pmatrix}$  的特征值, 因此

$$(\alpha_1 - \sigma)(\alpha_2 - \sigma) = \beta_1^2,$$

从 (33) 可知  $|\alpha_1 - \sigma| \leq \beta_1$ ,

等式只有当  $\delta = 0$  时才成立. 同样

$$|\alpha_1 - \sigma| \leq |\alpha_2 - \sigma|,$$

等式只有当  $\delta = 0$  时成立. 这两个关系, 可以写成

$$\frac{|\alpha_1 - \sigma|}{|\beta_1|} = \frac{|\beta_1|}{|\alpha_2 - \sigma|} = \sqrt{\frac{|\alpha_1 - \sigma|}{|\alpha_2 - \sigma|}} \leq 1,$$

最后的  $\leq$  号中的等号只有当  $\delta = 0$  时才成立.

现在来考察, 带 Wilkinson 位移时,  $\det(T - \sigma I)$  是什么?

从

$$\det(T - \sigma I) = (\alpha_1 - \sigma) \det(T_{2,n} - \sigma I) - \beta_1^2 \det(T_{3,n} - \sigma I),$$

而

$$\begin{aligned} \det(T_{2,n} - \sigma I) &= (\alpha_2 - \sigma) \det(T_{3,n} - \sigma I) \\ &\quad - \beta_2^2 \det(T_{4,n} - \sigma I), \end{aligned}$$

于是

$$\begin{aligned} \det(T - \sigma I) &= (\alpha_1 - \sigma) [(\alpha_2 - \sigma) \det(T_{3,n} - \sigma I) \\ &\quad - \beta_2^2 \det(T_{4,n} - \sigma I)] - \beta_1^2 \det(T_{3,n} - \sigma I) \\ &= [(\alpha_1 - \sigma)(\alpha_2 - \sigma) - \beta_1^2] \det(T_{3,n} - \sigma I) \\ &\quad - (\alpha_1 - \sigma) \beta_2^2 \det(T_{4,n} - \sigma I), \end{aligned}$$

因为  $(\alpha_1 - \sigma)(\alpha_2 - \sigma) - \beta_1^2 = 0$ , 故

$$\det(T - \sigma I) = -(\alpha_1 - \sigma) \beta_2^2 \det(T_{4,n} - \sigma I). \quad (34)$$

**定理 3.13** 设对任意对称不可约三对角阵  $T$ , 使用带 Wilkinson 位移  $\{\sigma_k\}$  的 QL 方法, 获矩阵序列  $\{T^{(k)}\}$ ,

$$T^{(k)} = \begin{pmatrix} \alpha_1^{(k)} & \beta_1^{(k)} & & & 0 \\ \beta_1^{(k)} & \alpha_2^{(k)} & \beta_2^{(k)} & & \\ & \beta_2^{(k)} & \ddots & \ddots & \\ & & \ddots & \ddots & \beta_{n-1}^{(k)} \\ 0 & & & \beta_{n-1}^{(k)} & \alpha_n^{(k)} \end{pmatrix},$$

则  $\beta_1^{(k)} \rightarrow 0$ .

证明 若某个  $\sigma_k$  是  $T$  的特征值, 则  $\beta_1^{(k+1)} = 0$  且  $\beta_1^{(n)} = 0$ ,  $l > k+1$ , 因此可设  $\sigma_k$  不是  $T$  的特征值. 为了简明起见, 在证明中省略序列指标  $k$ , 以  $\hat{\alpha}_i, \hat{\beta}_i$  表示  $T^{(k+1)}$  的元素, 以  $\alpha_i, \beta_i$  表示  $T^{(k)}$  的元素.

利用 § 2 定理 3.9, 对一切位移成立

$$l_{11}^2 = \det(T - \sigma I)^2 / \omega_1,$$

$$\omega_1 = \det(T_{2,n} - \sigma I)^2 + (\det(T_{3,n} - \sigma I) \beta_1)^2 \\ + (\det(T_{4,n} - \sigma I) \beta_1 \beta_2)^2 + \cdots + (\beta_1 \beta_2 \cdots \beta_{n-1})^2,$$

$$(T - \sigma I) q_1 = l_{11} e_1,$$

$$q_{11} = l_{11} \det(T_{2,n} - \sigma I) / \det(T - \sigma I),$$

$$q_{21} = -l_{11} \beta_1 \det(T_{3,n} - \sigma I) / \det(T - \sigma I),$$

$$\text{满足} \quad (\alpha_1 - \sigma) q_{11} + \beta_1 q_{21} = l_{11},$$

或者

$$(\alpha_1 - \sigma) \frac{\det(T_{2,n} - \sigma I)}{\det(T - \sigma I)} + \beta_1 \left( -\beta_1 \frac{\det(T_{3,n} - \sigma I)}{\det(T - \sigma I)} \right) = 1,$$

因此点  $\left( \frac{\det(T_{2,n} - \sigma I)}{\det(T - \sigma I)}, -\beta_1 \frac{\det(T_{3,n} - \sigma I)}{\det(T - \sigma I)} \right)$  在直线

$$(\alpha_1 - \sigma)x + \beta_1 y = 1$$

上. 而该直线离原点的最短距离为  $1/\sqrt{(\alpha_1 - \sigma)^2 + \beta_1^2}$ , 故

$$\begin{aligned}
& \left( \frac{\det(T_{2,n}-\sigma I)}{\det(T-\sigma I)} \right)^2 + \left( \beta_1 \frac{\det(T_{3,n}-\sigma I)}{\det(T-\sigma I)} \right)^2 \\
& \geq \frac{1}{(\alpha_1-\sigma)^2 + \beta_1^2}, \\
l_{11}^2 & \leq 1 / \left[ \left( \frac{\det(T_{2,n}-\sigma I)}{\det(T-\sigma I)} \right)^2 + \left( \beta_1 \frac{\det(T_{3,n}-\sigma I)}{\det(T-\sigma I)} \right)^2 \right] \\
& \leq 1 / ((\alpha_1-\sigma)^2 + \beta_1^2)^{-1} = (\alpha_1-\sigma)^2 + \beta_1^2, \quad (35)
\end{aligned}$$

这是对一切位移 $\sigma$ 都成立的估计式。现在假如取的 $\sigma$ 是Wilkinson位移,于是

$$(\alpha_1 - \sigma)^2 \leq \beta_1^2,$$

故

$$l_{11}^2 \leq 2\beta_1^2. \quad (36)$$

再从

$$\begin{aligned}
l_{11}^2 & \leq 1 / \left[ \left( \frac{\det(T_{2,n}-\sigma I)}{\det(T-\sigma I)} \right)^2 \right. \\
& \quad \left. + \left( \frac{\beta_1 \det(T_{3,n}-\sigma)}{\det(T-\sigma)} \right)^2 + \left( \beta_1 \beta_2 \frac{\det(T_{4,n}-\sigma)}{\det(T-\sigma)} \right)^2 \right].
\end{aligned}$$

利用

$$\det(T-\sigma I) = -(\alpha_1 - \sigma) \beta_2^2 \det(T_{4,n} - \sigma I),$$

故

$$\begin{aligned}
l_{11}^2 & \leq 1 / \left[ \frac{1}{(\alpha_1 - \sigma)^2 + \beta_1^2} + \left( \frac{\beta_1 \beta_2}{(\alpha_1 - \sigma) \beta_2^2} \right)^2 \right] \\
& \leq 1 / \left( \frac{1}{2\beta_1^2} + \left( \frac{\beta_1}{(\alpha_1 - \sigma) \beta_2} \right)^2 \right) \\
& \leq 1 / \left( \frac{1}{2\beta_1^2} + \left( \frac{1}{\beta_2} \right)^2 \right) = \frac{2\beta_1^2 \beta_2^2}{\beta_2^2 + 2\beta_1^2} \\
& = \frac{(\sqrt{2} |\beta_1|) |\beta_2| (\sqrt{2} |\beta_1|) |\beta_2|}{(\sqrt{2} \beta_1)^2 + \beta_2^2} \leq \frac{1}{2} \sqrt{2} |\beta_1| |\beta_2|. \quad (37)
\end{aligned}$$

下面先证  $\beta_1^2 \beta_2 \rightarrow 0$ , 因为

$$\hat{T} = Q^* T Q,$$

是从一个对称三对角阵正交相似变换到另一对称三对角阵, 故可应用第2章 §3 的极值性质

$$|\hat{\beta}_1 \hat{\beta}_2| = \min_{\psi_2 \in m p_2} \|\psi_2(T) q_1\|,$$

特别取  $\psi_2(\lambda) = (\lambda - \alpha_1)(\lambda - \sigma)$ , 于是

$$\begin{aligned} |\hat{\beta}_1 \hat{\beta}_2| &\leq \| (T - \alpha_1)(T - \sigma) q_1 \| \\ &= \| (T - \alpha_1) l_{11} e_1 \| = l_{11} |\beta_1|, \end{aligned}$$

利用定理 3.8 的推论  $|\hat{\beta}_1| \leq l_{11}$ , 故

$$|\hat{\beta}_1^2 \hat{\beta}_2| \leq l_{11}^2 |\beta_1|,$$

将(37)代入上式, 即得

$$|\hat{\beta}_1^2 \hat{\beta}_2| \leq \frac{\sqrt{2}}{2} |\beta_1^2 \beta_2|,$$

而  $\sqrt{2}/2 = 0.7071067 \dots < 1$ , 这就证明了

$$\beta_1^2 \beta_2 \rightarrow 0.$$

另一方面

$$\begin{aligned} |\hat{\beta}_1|^3 &= |\hat{\beta}_1| |\hat{\beta}_1|^2 \leq l_{11} l_{11}^2 \\ &\leq \sqrt{2} |\beta_1| \frac{\sqrt{2}}{2} |\beta_1| |\beta_2| = |\beta_1^2 \beta_2|, \end{aligned}$$

故  $|\beta_1^3| \rightarrow 0$ ,

从而证明了  $\beta_1 \rightarrow 0$ . 证毕.

下面的定理给出收敛速度的估计.

**定理 3.14** 对任意对称不可约三对角阵  $T$ , 使用带 Wilkinson 位移的 QL 方法, 则  $T^{(k+1)}$  的  $\hat{\beta}_1$  与  $T^{(k)}$  的  $\beta_1$  之间有下列关系:

$$|\hat{\beta}_1| = |\alpha_1 - \sigma_k| |\beta_1| \beta_2^2 |\det(T_{4,n}^{(k)} - \sigma_k I)| (\omega_2^{(k)})^{\frac{1}{2}} / |\omega_1^{(k)}|, \quad (38)$$

这里

$$\begin{aligned}\omega_1^{(k)} &= (\det(T_{2,n}^{(k)} - \sigma_k I))^2 + (\det(T_{3,n}^{(k)} - \sigma_k I) \beta_1)^2 \\ &\quad + (\det(T_{4,n}^{(k)} - \sigma_k I) \beta_1 \beta_2)^2 + \cdots + (\beta_1 \beta_2 \cdots \beta_{n-1})^2 \\ &\rightarrow \left( \prod_{\substack{l=1 \\ l \neq 4}}^n (\theta_l - \theta_i) \right)^2,\end{aligned}$$

其中  $\theta_i$  是  $\alpha_1$  的极限.

$$\begin{aligned}\omega_2^{(k)} &= (\det(T_{3,n}^{(k)} - \sigma_k I))^2 + (\det(T_{4,n}^{(k)} - \sigma_k I) \beta_2)^2 \\ &\quad + (\det(T_{5,n}^{(k)} - \sigma_k I) \beta_2 \beta_3)^2 + \cdots + (\beta_2 \beta_3 \cdots \beta_{n-1})^2\end{aligned}$$

是一致有界量.

特别当

$$\lim_{k \rightarrow \infty} |\alpha_2^{(k)} - \sigma_k| \neq 0, \quad (39)$$

有关系式:

$$|\hat{\beta}_1| = |\beta_1|^3 \beta_2^2 |\det(T_{4,n}^{(k)} - \sigma_k I)| (\omega_2^{(k)})^{\frac{1}{2}} / (|\alpha_2 - \sigma_k| \omega_1^{(k)}). \quad (40)$$

证明 因为  $\beta_1 \rightarrow 0$ , 故  $\alpha_1 - \sigma_k \rightarrow 0$ , 因此  $\alpha_1$  的极限存在, 且是  $T$  的特征值  $\theta_i$ , 参见 [31].

利用定理 3.10 有

$$\hat{\beta}_1^2 = \beta_1^2 \det(T^{(k)} - \sigma_k I)^2 \omega_2^{(k)} / (\omega_1^{(k)})^2,$$

类似于定理 3.12 对于带 Rayleigh 位移的讨论, 利用  $\beta_1 \rightarrow 0$  有

$$\omega_1^{(k)} \rightarrow \prod_{\substack{l=1 \\ l \neq i}}^n (\theta_l - \theta_i)^2$$

和  $\omega_2^{(k)}$  是一致有界量.

再利用  $\sigma_k$  是 Wilkinson 位移, 故成立

$$\det(T^{(k)} - \sigma_k I) = -(\alpha_1 - \sigma_k) \beta_2^2 \det(T_{4,n}^{(k)} - \sigma_k I),$$

于是  $\hat{\beta}_1^2 = (\alpha_1 - \sigma_k)^2 \beta_1^2 \beta_2^4 \det(T_{4,n}^{(k)} - \sigma_k I)^2 \omega_2^{(k)} / (\omega_1^{(k)})^2$ ,

即 (38) 式成立.

特别当  $\lim_{k \rightarrow \infty} |\alpha_2^{(k)} - \sigma_k| \neq 0$ , 利用  $\alpha_1 - \sigma_k = \beta_1^2 / (\alpha_2 - \sigma_k)$  故得

$$|\hat{\beta}_1| = |\beta_1|^3 |\beta_2^2| \det(T_{4,n}^{(k)} - \sigma_k I) (\omega_2^{(k)})^{\frac{1}{2}} / (|\alpha_2 - \sigma_k| \omega_1^{(k)}),$$
 即(40)式成立. 证毕.

从(38)式利用  $|\alpha_1 - \sigma_k| \leq \beta_1$  可以知道,  $\beta_1$  的收敛速度至少是二次的. 当条件(39)成立时, 从(40)式可知收敛速度至少是三次的. 因为  $\sigma_k \rightarrow \theta_i$ , 如果  $\beta_3 \rightarrow 0$ , 那么  $\alpha_2 \rightarrow \theta_i \neq \theta_i$ , 此时(39)一定成立, 并且  $|\hat{\beta}_1| = O(|\beta_1|^3 |\beta_2^2|)$  成立. 于是可知要  $|\hat{\beta}_1| = O(|\beta_1|^3 |\beta_2^2|)$  用不着再附加条件  $\beta_3 \rightarrow 0$ . 这一点已说明定理 3.14 要比[10]中结果更强一些.

关于其他一些位移下, QL 方法的收敛性和收敛率的讨论, 可参看[31]和[32]. 在[32]中提出了一种称为 RW 位移的新的位移, 它也保证对任意不可约对称三对角矩阵  $T$  情况下的  $\beta_1 \rightarrow 0$ , 并且在任何情况下  $\beta_1$  收敛率至少是三次的.

## 解特征值问题的 Lanczos 算法

Lanczos 算法求对称矩阵的特征值是 1950 年提出的 [23], 取定一个单位向量  $q$ , 通过第 2 章所介绍 Lanczos 过程 (7) 构造一组正交化序列  $q_1, q_2, \dots, q_n$ ,  $Q = [q_1, q_2, \dots, q_n]$ , 则  $Q^*AQ = T$  是一个对称三对角矩阵. 这样矩阵  $A$  的特征值问题, 化成矩阵  $T$  的特征值问题, 后者自然比较简单多了. 当然  $q_1, q_2$  计算到  $q_j (j < n)$ ,  $Q_j = [q_1, q_2, \dots, q_j]$  得到  $j \times j$  对称三对角阵  $T_j = Q_j^*AQ_j$ , 可以求  $T_j$  的特征值, 作为  $A$  的特征值近似值, 也是把求  $A$  的特征值问题化成简单的求  $T_j$  的特征值问题了. 这就是 Lanczos 算法求特征值的基本思想.

但是由于舍入误差的影响, 在求  $q_i$  的过程中, 所求到的  $q_i$  很快就失去了正交性, 因此 Lanczos 算法 50 年代、60 年代被认为是不稳定的, 很少被人使用. 直到 1971 年 O. C. Paige 在他的著名博士论文 [13] 中, 通过舍入误差的分析, 发现失去正交性, 恰与近似特征值的精度提高相关. 重新肯定 Lanczos 算法是一种求大稀疏矩阵, 两端部的特征值, 即最大的部分和最小的部分特征值的一种有效方法. 自那以后 Lanczos 算法越来越被重视, 这是因为它具有很大的优点:

1. 每次迭代用到原来矩阵  $A$  的只是  $A$  与向量的乘法, 因此可以充分利用  $A$  的稀疏性.
2. 在一次迭代中要进行的计算, 都是向量间的运算, 有利于实现平行运算.
3. 在一次迭代中, 只用到二个向量, 其余已算得的资料



可以放在外存贮中.

4. 对计算的结果, 可以有误差的估计.

因为有上述这些优点, 因此目前对于非常高阶的稀疏对称矩阵的特征值问题, 都用 Lanczos 算法来解.

## § 1 近似不变子空间

设  $A$  是  $n \times n$  实对称矩阵, 它的  $n$  个特征值, 由小到大排成  $\lambda_1, \lambda_2, \dots, \lambda_n$ , 对应的单位正交特征向量为  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$ .

对任意非零  $\mathbf{x} \in R^n$ , 定义 Rayleigh 商为

$$R(\mathbf{x}) = (A\mathbf{x}, \mathbf{x}) / (\mathbf{x}, \mathbf{x})$$

有如下极大极小原理

**定理 4.1** 记  $V_l$  是一个任意  $l$  维子空间,  $l \leq n$ , 则

$$\lambda_l = \min_{V_l} \max_{\mathbf{x} \in V_l} R(\mathbf{x}), \quad (1)$$

$$\lambda_{n-l+1} = \max_{V_l} \min_{\mathbf{x} \in V_l} R(\mathbf{x}). \quad (2)$$

**证明** 记  $n-l+1$  维空间  $E_{l-1}^\perp = \{\mathbf{y}_l, \mathbf{y}_{l+1}, \dots, \mathbf{y}_n\}$ , 于是任意  $l$  维空间  $V_l$  与  $E_{l-1}^\perp$  的交不空. 取非零向量  $\mathbf{y} \in V_l \cap E_{l-1}^\perp$  有

$$\max_{\mathbf{x} \in V_l} R(\mathbf{x}) \geq R(\mathbf{y}) \geq \lambda_l,$$

故

$$\min_{V_l} \max_{\mathbf{x} \in V_l} R(\mathbf{x}) \geq \lambda_l.$$

另一方面, 取一个  $l$  维子空间  $\hat{V}_l = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_l\}$ ,

$$\max_{\mathbf{x} \in \hat{V}_l} R(\mathbf{x}) = \lambda_l \geq \min_{V_l} \max_{\mathbf{x} \in V_l} R(\mathbf{x}),$$

从而(1)成立. 同样可证(2)成立. 证毕.

**定理 4.2** 设  $E$  也是  $n \times n$  对称矩阵, 它的特征值由小到大排成  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ , 如果将  $A+E$  的  $n$  个特征值由小到大

排成  $\mu_1, \mu_2, \dots, \mu_n$ , 则

$$\mu_l \leq \lambda_l + \varepsilon_n, \quad (3)$$

$$\lambda_l \leq \mu_l - \varepsilon_1. \quad (4)$$

证明 由定理 4.1

$$\begin{aligned} \mu_l &= \min_{V_l} \max_{x \in V_l} ((A+E)x, x) / (x, x) \\ &= \min_{V_l} \max_{x \in V_l} \left[ \frac{(Ax, x)}{(x, x)} + \frac{(Ex, x)}{(x, x)} \right] \\ &\leq \max_{x \in \hat{V}_l} \left[ \frac{(Ax, x)}{(x, x)} + \frac{(Ex, x)}{(x, x)} \right], \end{aligned}$$

这里  $\hat{V}_l = \{y_1, y_2, \dots, y_l\}$ . 于是

$$\mu_l \leq \max_{x \in \hat{V}_l} \frac{(Ax, x)}{(x, x)} + \max_{x \in \hat{V}_l} \frac{(Ex, x)}{(x, x)} \leq \lambda_l + \varepsilon_n.$$

由  $A = A + E - E$ , 而  $-E$  的最大特征值为  $-\varepsilon_1$ , 利用上述结果就得

$$\lambda_l \leq \mu_l - \varepsilon_1. \text{ 证毕.}$$

推论

$$|\lambda_l - \mu_l| \leq \|E\|. \quad (5)$$

设  $Q = [q_1, q_2, \dots, q_j]$  是  $j$  个  $n$  维、线性无关的列向量组成的矩阵,  $j < n$ , 若存在  $j \times j$  矩阵  $T_j$  使

$$AQ - QT_j = 0,$$

那么  $\{q_1, q_2, \dots, q_j\}$  是  $A$  的不变子空间. 此时  $T_j$  的特征值, 即为  $A$  的特征值, 若  $s$  是  $T_j$  的特征向量, 则  $Qs$  是  $A$  的特征向量. 如果对于  $Q$ , 有  $T_j$  使得

$$AQ - QT_j = R,$$

$\|R\|$  很小, 那么可以设想  $\{q_1, q_2, \dots, q_n\}$  是  $A$  的近似不变子空间. 下面我们要讨论此时  $T_j$  的特征值、特征向量与  $A$  的特征值、特征向量之间的关系.

对于给定的  $Q$ , 任意取一个  $j \times j$  矩阵  $T_j$ , 都有一个  $R(T_j)$ , 使得

$$AQ - QT_j = R(T_j),$$

首先有一个问题是: 选取怎样的  $T_j$ , 使  $\|R(T_j)\|$  最小.

对于  $j=1$  的情况,  $T_1$  是个常数  $\alpha$ , 考虑

$$Aq_1 - q_1\alpha = r,$$

问  $\alpha$  取什么值时, 使  $\|r\| = \|Aq_1 - q_1\alpha\|$  最小. 这问题可以看成向量  $Aq_1$ , 在  $q_1$  所张成的子空间中找一个最佳逼近向量, 利用第 1 章的定理 1.1, 我们知道当  $\alpha$  使

$$(Aq_1 - \alpha q_1, q_1) = 0$$

时,  $\|r\|$  最小, 也即取

$$\alpha^0 = (Aq_1, q_1) / (q_1, q_1), \quad (6)$$

有  $\|Aq_1 - \alpha^0 q_1\| \leq \|Aq_1 - \alpha q_1\|, \quad \forall \alpha \in R,$

注意(6)式恰为  $q_1$  的 Rayleigh 商, 因此可以说: 对于  $Aq_1 - q_1\alpha$  当  $\alpha$  取成  $q_1$  的 Rayleigh 商时, 使  $\|Aq_1 - q_1\alpha\|$  最小.

对于任意的  $\alpha$ , 它跟  $A$  的特征值之间成立如下定理中所述的估计式.

**定理 4.3** 对于  $\alpha$ , 可以找到一个  $A$  的特征值  $\lambda$ , 使得

$$|\lambda - \alpha| \leq \|r\| / \|q_1\|. \quad (7)$$

**证明** 将  $q_1$  按  $A$  的特征向量展开

$$q_1 = \sum_{i=1}^n \beta_i y_i,$$

于是

$$r = \sum_{i=1}^n \beta_i (\lambda_i - \alpha) y_i,$$

$$\|r\|^2 = \sum_{i=1}^n \beta_i^2 (\lambda_i - \alpha)^2$$

$$\geq \min_i (\lambda_i - \alpha)^2 \sum_{i=1}^n \beta_i^2 = \min_i (\lambda_i - \alpha)^2 \|q_1\|^2,$$

故  $\min_i (\lambda_i - \alpha)^2 \leq \|r\|^2 / \|q_1\|^2$ . 证毕.

在定理证明中使  $(\lambda_i - \alpha)^2$  达到极小的那个特征值  $\lambda$ , 如果对应的特征向量不止一个, 那末  $q_1$  的分解式, 可以写成

$$q_1 = \sum_{\lambda_i = \lambda} \beta_i y_i + \sum_{\lambda_i \neq \lambda} \beta_i y_i,$$

将  $\sum_{\lambda_i = \lambda} \beta_i y_i$  重新写成  $\beta y$ , 其中  $y$  是单位向量, 显然也是对应特征值  $\lambda$  的特征向量. 于是

$$q_1 = \beta y + \sum_{\lambda_i \neq \lambda} \beta_i y_i,$$

假定  $\|q_1\| = 1$ , 可得  $\beta^2 + \sum_{\lambda_i \neq \lambda} \beta_i^2 = 1$ .

**定理 4.4** 如果  $d = \min_{\lambda_i \neq \lambda} |\lambda_i - \alpha|$ , 则

$$\|q_1 - \beta y\| \leq \|r\| / d.$$

**证明** 从  $Aq_1 - \alpha q_1 = r$  有

$$\beta(\lambda - \alpha)y + \sum_{\lambda_i \neq \lambda} \beta_i(\lambda_i - \alpha)y_i = r,$$

$$\|r\|^2 \geq \sum_{\lambda_i \neq \lambda} \beta_i^2 (\lambda_i - \alpha)^2 \geq d^2 \sum_{\lambda_i \neq \lambda} \beta_i^2,$$

故  $\sum_{\lambda_i \neq \lambda} \beta_i^2 \leq \|r\|^2 / d^2$ ,

但  $q_1 - \beta y = \sum_{\lambda_i \neq \lambda} \beta_i y_i$ ,

故  $\|q_1 - \beta y\|^2 \leq \sum_{\lambda_i \neq \lambda} \beta_i^2 \leq \|r\|^2 / d^2$ . 证毕.

**推论** 若  $q_1$  与  $y$  的交角  $\angle(q_1, y) = \theta$ , 则

$$|\sin \theta| \leq \|r\| / d. \quad (8)$$

**证明** 因为  $(q_1, y) = \beta = \cos \theta$ , 故  $\sin^2 \theta = 1 - \beta^2$ ,

$$\sin^2 \theta = \sum_{\lambda_i \neq \lambda} \beta_i^2,$$

从而  $|\sin \theta| \leq \|r\| / d$ .

上述定理描写了  $q_1$  逼近  $A$  的特征向量的状况. 特别当  $\alpha$  取成  $q_1$  的 Rayleigh 商  $\rho(q_1)$  时, 有

**定理 4.5** 设  $\|q_1\|=1$ ,  $\rho=\rho(q_1)$ ,  $r=Aq_1-\rho q_1$ ,  $d=\min_{\lambda_i \neq \lambda} |\lambda_i - \rho|$ ,  $\lambda$  是与  $\rho$  最接近的特征值, 则

$$|\lambda - \rho| \leq \|r\|^2/d. \quad (9)$$

**证明** 利用分解式

$$q_1 = \beta y + \sum_{\lambda_i \neq \lambda} \beta_i y_i,$$

$$r = \beta(\lambda - \rho)y + \sum_{\lambda_i \neq \lambda} \beta_i(\lambda_i - \rho)y_i,$$

有  $\|r\|^2 = \beta^2(\lambda - \rho)^2 + \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho)^2.$

利用  $q_1^* r = 0$ , 故

$$\beta^2(\lambda - \rho) + \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho) = 0,$$

或  $\beta^2(\lambda - \rho) = - \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho),$

代入  $\|r\|^2$  的表示式中

$$\begin{aligned} \|r\|^2 &= -(\lambda - \rho) \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho) + \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho)^2 \\ &= \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \lambda)(\lambda_i - \rho), \end{aligned}$$

又从  $\rho - \lambda = q_1^*(A - \lambda)q_1 = \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \lambda),$

故

$$|\rho - \lambda| \leq \sum_{\lambda_i \neq \lambda} \beta_i^2 |\lambda_i - \lambda| \leq \frac{1}{d} \sum_{\lambda_i \neq \lambda} \beta_i^2 |\lambda_i - \lambda| |\lambda_i - \rho|. \quad (10)$$

因为  $\lambda$  是  $\rho$  最接近的特征值, 故有

$$|\lambda_i - \lambda| |\lambda_i - \rho| = (\lambda_i - \lambda)(\lambda_i - \rho),$$

将此代入(10)即得

$$|\rho - \lambda| \leq \frac{1}{d} \sum_{\lambda_i \neq \lambda} \beta_i^2 (\lambda_i - \lambda)(\lambda_i - \rho) = \|r\|^2/d. \text{ 证毕.}$$

**推论** 若  $\lambda$  是  $\lambda_1$  或  $\lambda_n$ , 则

$$|\tan \theta| \leq \|r\|/d.$$

证明 由

$$\beta^2(\lambda - \rho) = - \sum_{\lambda_i \neq \lambda} \beta_i^2(\lambda_i - \rho),$$

因为所有  $\lambda_i \neq \lambda$  的  $\lambda_i - \rho$  都是同号的, 因此

$$|\lambda - \rho| \beta^2 = \sum_{\lambda_i \neq \lambda} \beta_i^2 |\lambda_i - \rho| \geq d \sum_{\lambda_i \neq \lambda} \beta_i^2,$$

$$\tan^2 \theta = \sum_{\lambda_i \neq \lambda} \beta_i^2 / \beta^2 \leq |\lambda - \rho| / d \leq \|r\|^2 / d^2. \text{ 证毕.}$$

这个定理和推论说明当  $\alpha$  取成 Rayleigh 商, 其估计有所加强.

现在来考虑,  $Q$  有  $j$  个列,  $j > 1$  的情况, 设  $Q = [q_1, q_2, \dots, q_j]$ ,  $q_1, q_2, \dots, q_j$  是单位正交向量组,

$$AQ - QT_j = R,$$

讨论  $T_j$  的特征值与  $A$  的特征值的关系. 首先有

**定理 4.6** 对任意  $T_j$  的特征值  $\mu$ , 可以找到一个  $A$  的特征值  $\lambda$ , 使得

$$|\lambda - \mu| \leq \|R\|.$$

**证明**  $\mu$  是  $T_j$  的特征值, 记  $s$  是它对应的单位特征向量, 令  $x = Qs$ , 于是

$$\|x\| = \|s^* Q^* Q s\|^{1/2} = 1,$$

$$AQs - QT_j s = Rs,$$

$$Ax - \mu x = Rs.$$

利用定理 4.3, 有  $A$  的特征值  $\lambda$  使

$$|\lambda - \mu| \leq \|Rs\| \leq \|R\|. \text{ 证毕.}$$

这个定理的缺点是: 对  $T_j$  的两个特征值  $\mu_1, \mu_2$ , 定理中所说的  $A$  的特征值  $\lambda$ , 可能是同一个. 即

$$|\lambda - \mu_1| \leq \|R\|,$$

$$|\lambda - \mu_2| \leq \|R\|.$$

人们希望对于  $T_j$  的  $j$  个特征值  $\mu_1, \mu_2, \dots, \mu_j$ , 对应找

到  $A$  的  $j$  个特征值  $\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_j}$ ,  $i_1, i_2, \dots, i_j$  是  $j$  个不同的数, 得到  $|\mu_k - \lambda_{i_k}|$  的估计式.

这个任务在 1967 年, W. Kahan 的没有公开发表的著名论文: "Inclusion theorems for clusters of eigenvalues of Hermitian matrices" 中得到完成. 在那篇论文中, 他先证明了一个重要定理, 可称为矩阵的保范扩张定理. 这个定理第一次公开发表在 B. N. Parlett 的书 [10] 中, 在 [10] 中的证明, B. N. Parlett 作了适当的修改. 下面我们给出这个矩阵保范扩张定理, 是采用 W. Kahan 原来的证明.

**定理 4.7** 设  $V$  是  $n \times n$  Hermite 矩阵,  $C$  是任意复数域上的  $m \times n$  矩阵, 对于  $(n+m) \times n$  矩阵

$$R = \begin{pmatrix} V \\ C \end{pmatrix}$$

存在  $m \times m$  Hermite 矩阵  $W$ , 使得  $(n+m) \times (n+m)$  Hermite 矩阵

$$T = \begin{pmatrix} V & C^* \\ C & W \end{pmatrix}$$

有  $\|T\| = \|R\|$ .

证明

$$R^*R = V^2 + C^*C,$$

$$T^2 = \begin{pmatrix} V^2 + C^*C & VC^* + C^*W \\ CV + WC & CC^* + W^2 \end{pmatrix},$$

$V^2 + C^*C$  是  $T^2$  的主子阵, 因此对任意 Hermite 矩阵  $W$ , 都有

$$\|T^2\| \geq \|R\|^2, \quad \text{及} \quad \|T\| \geq \|R\|. \quad (11)$$

下面要找一个  $W$ , 使得  $\|T\|$  最小. 取任意数  $\rho > \|R\|$ , 于是  $\rho^2 I - V^2$  是一个 Hermite 正定矩阵, 当然有逆. 令  $D = C(\rho^2 I - V^2)^{-1}$ ,

$$\rho^2 I - T^2 = \begin{pmatrix} \rho^2 I - V^2 - C^* C & -VC^* - C^* W \\ -CV - WC & \rho^2 I - CC^* - W^2 \end{pmatrix},$$

考察

$$\begin{pmatrix} I & 0 \\ DV & I \end{pmatrix} (\rho^2 I - T^2) \begin{pmatrix} I & VD^* \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & 0 \\ DV & I \end{pmatrix} \\ \times \begin{pmatrix} \rho^2 I - V^2 - C^* C & (\rho^2 I - V^2 - C^* C) VD^* - VC^* - C^* W \\ * & (-CV - WC) VD^* + \rho^2 I - CC^* - W^2 \end{pmatrix},$$

因此如果取  $W$  使

$$\begin{aligned} & (\rho^2 I - V^2 - C^* C) VD^* - VC^* - C^* W \\ & = -C^* C VD^* - C^* W = 0, \end{aligned}$$

上述三个矩阵相乘的积, 就是一个块对角阵. 为此取

$$W(\rho) = -CV D^* = -DVC^*,$$

此时

$$\begin{pmatrix} I & 0 \\ DV & I \end{pmatrix} (\rho^2 I - T^2) \begin{pmatrix} I & VD^* \\ 0 & I \end{pmatrix} = \begin{pmatrix} \rho^2 I - R^* R & 0 \\ 0 & \rho^2 X \end{pmatrix}, \\ X = I - C(\rho^2 I - V^2)^{-1} C^*.$$

另一方面

$$\begin{aligned} & \begin{pmatrix} I & 0 \\ DV & I \end{pmatrix} (\rho^2 I - RR^*) \begin{pmatrix} I & VD^* \\ 0 & I \end{pmatrix} \\ & = \begin{pmatrix} I & 0 \\ DV & I \end{pmatrix} \begin{pmatrix} \rho^2 I - V^2 & -VC^* \\ -CV & \rho^2 I - CC^* \end{pmatrix} \begin{pmatrix} I & DV \\ 0 & I \end{pmatrix} \\ & = \begin{pmatrix} \rho^2 I - V^2 & 0 \\ 0 & \rho^2 X \end{pmatrix}. \end{aligned}$$

因  $RR^*$  的最大特征值与  $R^*R$  的最大特征值相同, 因此矩阵  $\rho^2 I - RR^*$  也是正定矩阵, 从而由惯性定理 (或称 Sylvester 定理) [9, p. 335] 知道  $\rho^2 X$  也必须是正定矩阵, 由此知道  $\rho^2 I - T^2$  也必须是正定矩阵, 即得



$$\rho > \|T\|,$$

再看

$$W(\rho) = -CV(\rho^2 I - V^2)^{-1}C^*$$

是  $\rho$  的有理函数, 当  $\rho \rightarrow \|R\|$  时, 如果  $\|W(\rho)\|$  有界, 则  $\rho = \|R\|$  不是  $W(\rho)$  的奇点 (因为有理函数只可能有极点),  $CC^* + W^2$  也是  $T^2$  的主子阵, 因此  $T^2$  的最大特征值大于  $CC^* + W^2$  的最大特征值, 因为  $CC^*$  是非负定矩阵, 故

$$\|W(\rho)\| \leq \|T\| < \rho$$

有界, 故  $\rho \rightarrow \|R\|$ ,  $W(\rho)$  的极限  $W_+(\|R\|)$  存在. 并且对于这个  $W(\|R\|)$ , 有

$$\|T\| \leq \|R\|,$$

联合 (11) 式, 即有

$$\|T\| = \|R\|. \text{ 证毕.}$$

现在可以证明如下的 Kahan 定理.

**定理 4.8** 设  $Q$  是  $j$  个  $n$  维单位列正交向量组成,  $T_j$  是  $j \times j$  对称矩阵,

$$R = AQ - QT_j,$$

若  $T_j$  的  $j$  个特征值为  $\mu_1, \mu_2, \dots, \mu_j$ , 则可以找到  $j$  个不同的自然数  $i_1, i_2, \dots, i_j$  使得

$$|\lambda_{i_k} - \mu_k| \leq \|R\|. \quad (12)$$

**证明** 将  $Q$  扩充成  $n \times n$  正交矩阵  $G$ ,  $G = [Q, Q_1]$ , 矩阵

$$G^*AG = \begin{pmatrix} Q^*AQ & Q^*AQ_1 \\ Q_1^*AQ & Q_1^*AQ_1 \end{pmatrix}$$

与  $A$  有相同的特征值.

$$\begin{aligned} G^*R &= G^*AGG^*Q - G^*QT_j \\ &= \begin{pmatrix} Q^*AQ & Q^*AQ_1 \\ Q_1^*AQ & Q_1^*AQ_1 \end{pmatrix} \begin{pmatrix} I \\ 0 \end{pmatrix} - \begin{pmatrix} I \\ 0 \end{pmatrix} T_j \\ &= \begin{pmatrix} Q^*AQ - T_j \\ Q_1^*AQ \end{pmatrix}, \end{aligned}$$

因为  $G^*$  是正交阵, 故  $\|G^*R\| = \|R\|$ .

利用定理 4.7, 扩充  $G^*R$  成

$$\begin{pmatrix} Q^*AQ - T_j & Q^*AQ_1 \\ Q_1^*AQ & W \end{pmatrix} = E,$$

有  $W$  存在, 使得对称阵  $E$  的范数  $\|E\| = \|R\|$ .

考察 
$$G^*AG - E = \begin{pmatrix} T_j & 0 \\ 0 & Q_1^*AQ_1 - W \end{pmatrix}.$$

利用定理 4.2, 可知上述等式右端矩阵的第  $l$  个特征值, 与  $G^*AG$  的第  $l$  个特征值之差的绝对值小于等于  $\|E\|$ . 如果  $T_j$  的特征值  $\mu_1, \mu_2, \dots, \mu_j$ , 在矩阵

$$\begin{pmatrix} T_j & 0 \\ 0 & Q_1^*AQ_1 - W \end{pmatrix}$$

的特征值中, 按由小到大次序排列的编号为  $i_1, i_2, \dots, i_j$ , 于是

$$|\lambda_{i_k} - \mu_k| \leq \|E\| = \|R\|. \text{ 证毕.}$$

这个定理加强了定理 4.6 的结果.

对于  $Q = [q_1, q_2, \dots, q_j]$ ,  $q_1, q_2, \dots, q_j$  是  $j$  个单位正交列向量的情况, 在等式

$$AQ - QT_j = R$$

中, 也有一个使  $\|R\|$  最小的最佳  $T_j$  问题. 下面定理回答了这个问题, 见 [10].

**定理 4.9** 对于给定的  $j$  个单位正交列向量组成的  $n \times j$  矩阵  $Q$ , 对任意  $j \times j$  矩阵  $T_j$ , 记

$$R(T_j) = AQ - QT_j,$$

若记  $j \times j$  对称矩阵  $H = Q^*AQ$ , 则

$$\|R(H)\| \leq \|R(T_j)\|.$$

证明

$$\begin{aligned} R(T_j)^* R(T_j) &= (AQ - QT_j)^* (AQ - QT_j) \\ &= Q^* A^2 Q - T_j^* Q^* A Q - Q^* A Q T_j + T_j^* T_j, \\ &= Q^* A^2 Q - H^2 + (H - T_j)^* (H - T_j), \end{aligned}$$

但  $R(H)^* R(H) = Q^* A^2 Q - H^2,$

故  $R(T_j)^* R(T_j) = R(H)^* R(H) + (H - T_j)^* (H - T_j),$

又因为  $(H - T_j)^* (H - T_j)$  是一个非负矩阵, 于是

$$\|R(H)\| \leq \|R(T_j)\|. \text{ 证毕.}$$

W. Kahan 还考虑: 矩阵  $Q$  的  $j$  个列不正交的情况下,  $T_j$  的特征值, 与  $A$  的特征值之间的关系. 在上面曾经指出过的那篇文章中他给出了如下定理, 证明是按照[10]中所给出的方式.

**定理 4.10** 设  $Q$  是任意  $j$  个线性无关的列向量构成的  $n \times j$  矩阵,  $T_j$  是任意  $j \times j$  实对称矩阵, 对于

$$R = AQ - QT_j,$$

若  $Q$  的最小奇异值为  $\sigma_1$ ,  $T_j$  的特征值按由小到大次序排列是  $\mu_1, \mu_2, \dots, \mu_j$ , 则存在  $i_1, i_2, \dots, i_j$ ,  $j$  个不同自然数, 使得

$$|\lambda_{i_k} - \mu_k| \leq \sqrt{2} \|R\| / \sigma_1. \quad (13)$$

**证明** 将  $Q$  进行奇异值分解, 见[9, p. 468~470]. 存在  $n \times n$  正交阵  $P_1$  和  $j \times j$  正交阵  $P_2$ , 使得

$$P_1 Q P_2 = \begin{pmatrix} J \\ 0 \end{pmatrix},$$

$J$  为  $j \times j$  对角阵  $\text{diag}(\sigma_j, \sigma_{j-1}, \dots, \sigma_1)$ . 对于

$$R = AQ - QT_j,$$

有  $P_1 R P_2 = P_1 A P_1^* P_1 Q P_2 - P_1 Q P_2^* T_j P_2.$

记  $\hat{R} = P_1 R P_2$ ,  $\hat{A} = P_1 A P_1^*$ ,  $\hat{T}_j = P_2^* T_j P_2$ , 有

$\|\hat{R}\| = \|R\|$ ,  $\hat{A}$  相似于  $A$ ,  $\hat{T}_j$  相似于  $T_j$ ,

$$\hat{R} = \hat{A} \begin{pmatrix} J \\ 0 \end{pmatrix} - \begin{pmatrix} J \\ 0 \end{pmatrix} \hat{T}_j,$$

若 
$$\hat{A} = \begin{pmatrix} H & B^* \\ B & U \end{pmatrix} = \begin{pmatrix} \hat{T}_j & 0 \\ 0 & U - X \end{pmatrix} + \begin{pmatrix} H - \hat{T}_j & B^* \\ B & X \end{pmatrix},$$

于是从定理 4.2 可知存在  $i_1, i_2, \dots, i_j$  使得

$$|\lambda_{i_k} - \mu_k| \leq \left\| \begin{pmatrix} H - \hat{T}_j & B^* \\ B & X \end{pmatrix} \right\|,$$

利用定理 4.7, 有  $X$  使

$$\left\| \begin{pmatrix} H - \hat{T}_j & B^* \\ B & X \end{pmatrix} \right\| = \left\| \begin{pmatrix} H - \hat{T}_j \\ B \end{pmatrix} \right\|, \quad (14)$$

下面来估计上式右端, 易知

$$\hat{R} = \begin{pmatrix} HJ - J\hat{T}_j \\ BJ \end{pmatrix},$$

因为  $Q$  的列线性无关, 因此  $\sigma_1 > 0$ , 故  $J^{-1}$  存在,

$$\|J^{-1}\| = 1/\sigma_1,$$

$$\|B\| = \|BJJ^{-1}\| \leq \|BJ\|/\sigma_1,$$

而  $HJ - J\hat{T}_j = (H - \hat{T}_j)J + (\hat{T}_jJ - J\hat{T}_j),$

上式右端第二项是一个斜对称矩阵, 对任意向量  $x$ , 有

$$x^*(\hat{T}_jJ - J\hat{T}_j)x = 0,$$

若对称矩阵  $H - \hat{T}_j$  的模最大的特征值为  $\lambda$ , 对应的单位特征向量为  $y$ , 于是

$$y^*(HJ - J\hat{T}_j)y = \lambda y^*Jy,$$

即得  $\|HJ - J\hat{T}_j\| \geq |\lambda| \sigma_1 = \|H - \hat{T}_j\| \sigma_1.$

再回到(14)式, 记

$$M = \begin{pmatrix} H - \hat{T}_j \\ B \end{pmatrix},$$

$$M^*M = (H - \hat{T}_j, B^*) \begin{pmatrix} H - \hat{T}_j \\ B \end{pmatrix} = (H - \hat{T}_j)^2 + B^*B,$$

所以

$$\begin{aligned} \|M^*M\| &= \|(H - \hat{T}_j)^2 + B^*B\| \leq \|H - \hat{T}_j\|^2 + \|B\|^2 \\ &\leq (\|HJ - J\hat{T}_j\|^2 + \|BJ\|^2) / \sigma_1^2, \\ \hat{R}^*\hat{R} &= (HJ - J\hat{T}_j)^*(HJ - J\hat{T}_j) + (BJ)^*(BJ), \end{aligned}$$

右端两项都是非负矩阵, 故

$$\begin{aligned} \|\hat{R}\| &\geq \|HJ - J\hat{T}_j\|, \\ \|\hat{R}\| &\geq \|BJ\|, \end{aligned}$$

从而得到  $\|M^*M\| \leq 2\|\hat{R}\|^2 / \sigma_1^2$ . 证毕.

当  $Q$  的各列单位正交时, 按本定理结论, 此时  $\sigma_1 = 1$ , 从 (13) 得

$$|\lambda_{i_k} - \mu_k| \leq \sqrt{2} \|R\|,$$

与定理 4.8 比较, 多了一个  $\sqrt{2}$  因子, 因此可以设想  $\sqrt{2}$  可以改成 1, 但这个问题至今没有解决.

## § 2 Lanczos 算法

给定一个单位向量  $q_1$ , 构造序列

$$q_1, Aq_1, A^2q_1, \dots, A^j q_1, \quad (15)$$

$\{q_1, Aq_1, \dots, A^j q_1\}$  称为 Крылов 子空间. 在求线性方程组解时, 系数矩阵  $A$  是对称正定的情况, 共轭斜量法就是利用 Крылов 空间提供讯息来构造近似解的. 对于对称矩阵  $A$  的特征值问题. 很多求特征值的方法, 也是利用 Крылов 子空间提供的讯息. 例如, 1931 年 Крылов 利用它来计算  $q_1$  的最小多项式. 这一方法在  $n$  很小时是可以用的, 当  $n$  较大时, 因为稳定性上的原因, 完全不适用了. 乘幂法也是利用 Крылов 子空

间, 求特征值、特征向量的办法, 不过常常只能计算得个别特征值和特征向量, 并且实际计算中收敛速度常很慢。

Lanczos 算法, 也是利用 Крылов 子空间提供的讯息, 来求  $A$  的近似特征值和近似特征向量的方法。它也是将序列 (15) 正交化的办法。作法是取定一个单位向量  $q_1$ , 按照第 2 章 § 2 的 (7) 式, 求得  $q_2, q_3, \dots, q_j$  和  $\alpha_1, \alpha_2, \dots, \alpha_j; \beta_1, \beta_2, \dots, \beta_j$ 。得到  $Q_j = [q_1, q_2, \dots, q_j]$ ,

$$T_j = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \\ & & \ddots & \ddots & \beta_{j-1} \\ 0 & & & \beta_{j-1} & \alpha_j \end{pmatrix}.$$

计算  $T_j$  的特征值, 作为  $A$  的特征值的近似值, 若  $T_j$  的特征向量是  $s$ , 则  $Q_j s$  作为  $A$  的特征向量的近似。具体程式如下:

1. 取  $q_1, \|q_1\|=1, u_1 = Aq_1, 1 \rightarrow j$ ;
2.  $\alpha_j = q_j^* u_j$ ;
3.  $r_j = u_j - q_j \alpha_j$ ;
4.  $\beta_j = \|r_j\|$ ;
5. 如果  $\beta_j$  充分小, 则转入计算  $T_j$  的特征值和特征向量;

否则  $q_{j+1} = r_j / \beta_j$

6.  $u_{j+1} = Aq_{j+1} - \beta_j q_j, j+1 \rightarrow j$ , 转 2;

假如不考虑舍入误差, 计算是无限位小数进行的, 则当计算到第  $j$  步, 我们得到向量序列  $q_1, q_2, \dots, q_j$  和  $r_j$ , 这是  $j+1$  个正交向量组, 实际上对于  $j=1$  的情况只有  $q_1$  和  $r_1$  两个向量, 有

$$(r_1, q_1) = (Aq_1 - \alpha_1 q_1, q_1) = (Aq_1, q_1) - \alpha_1 (q_1, q_1),$$

按  $\alpha_1$  的定义知  $(r_1, q_1) = 0$ . 因此  $q_1, q_2, \dots, q_j$  和  $r_j$  是一组正交化向量组对  $j=1$  成立, 假定对  $j=k$  成立, 即  $q_1, q_2, \dots, q_k$  和  $r_k$  是一组正交化组, 假定  $r_k \neq 0$ . 考虑  $q_1, q_2, \dots, q_k, q_{k+1}$  和  $r_{k+1}$ . 因为  $r_k = \beta_k q_{k+1}$ , 故  $q_1, q_2, \dots, q_k, q_{k+1}$  是一组单位正交向量组, 为此只要证  $(r_{k+1}, q_l) = 0, l=1, 2, \dots, k, k+1$ , 对于  $l < k$ ,

$$\begin{aligned}(r_{k+1}, q_l) &= (Aq_{k+1} - \alpha_{k+1}q_{k+1} - \beta_k q_k, q_l) \\ &= (Aq_{k+1}, q_l) = (q_{k+1}, Aq_l) \\ &= (q_{k+1}, \beta_l q_{l+1} + \alpha_l q_l + \beta_{l-1} q_{l-1}) = 0.\end{aligned}$$

又

$$\begin{aligned}(r_{k+1}, q_k) &= (Aq_{k+1} - \alpha_{k+1}q_{k+1} - \beta_k q_k, q_k) \\ &= (Aq_{k+1}, q_k) - \beta_k (q_k, q_k) \\ &= (q_{k+1}, Aq_k) - \beta_k \\ &= (q_{k+1}, \beta_k q_{k+1} - \alpha_k q_k - \beta_{k-1} q_{k-1}) - \beta_k \\ &= \beta_k (q_{k+1}, q_{k+1}) - \beta_k = 0,\end{aligned}$$

而

$$\begin{aligned}(r_{k+1}, q_{k+1}) &= (Aq_{k+1} - \alpha_{k+1}q_{k+1} - \beta_k q_k, q_{k+1}) \\ &= (Aq_{k+1}, q_{k+1}) - \alpha_{k+1} (q_{k+1}, q_{k+1}) = 0.\end{aligned}$$

于是我们证明了命题: 若  $\beta_1, \beta_2, \dots, \beta_{j-1}$  都不为 0, 则  $q_1, q_2, \dots, q_j$  是单位正交向量组,  $q_1, q_2, \dots, q_j$  与  $r_j$  正交.

$Q_j = [q_1, q_2, \dots, q_j]$  有  $Q_j^* Q_j = I$ ,  $T_j$  是一个对称三对角不可约矩阵, 第一条超对角线上元素都为正的. 称  $T_j$  的特征值为  $A$  的 Ritz 值. 若  $\mu_i$  是  $T_j$  的特征值,  $s_i$  是对应  $\mu_i$  的单位特征向量, 称  $z_i = Q_j s_i$  为  $A$  的 Ritz 向量.

将  $r_j$  表示成  $\beta_j q_{j+1}$ , 有下列定理.

**定理 4.11**  $A, Q_j, T_j, \beta_j, q_{j+1}$  之间成立如下的关系

$$AQ_j - Q_j T_j = \beta_j \mathbf{q}_{j+1} \mathbf{e}_j^*. \quad (17)$$

证明 由构造过程

$$\beta_j \mathbf{q}_{j+1} = A \mathbf{q}_j - \alpha_j \mathbf{q}_j - \beta_{j-1} \mathbf{q}_{j-1},$$

得

$$A \mathbf{q}_j = \beta_j \mathbf{q}_{j+1} + \alpha_j \mathbf{q}_j + \beta_{j-1} \mathbf{q}_{j-1},$$

因此矩阵

$$AQ_j - Q_j T_j$$

除了第  $j$  列外, 其余各列都为 0. 但它的第  $j$  列为

$$A \mathbf{q}_j - \beta_{j-1} \mathbf{q}_{j-1} - \alpha_j \mathbf{q}_j = \beta_j \mathbf{q}_{j+1},$$

$n \times j$  矩阵  $\beta_j \mathbf{q}_{j+1} \mathbf{e}_j^*$  的第  $j$  列为  $\beta_j \mathbf{q}_{j+1}$ , 其余各列都为 0, 因此 (17) 成立. 证毕.

推论 存在  $i_1, i_2, \dots, i_j$   $j$  个不相同的自然数, 使得

$$|\lambda_{i_k} - \mu_k| \leq \beta_j.$$

证明 应用定理 4.8, 此时  $R = \beta_j \mathbf{q}_{j+1} \mathbf{e}_j^*$ , 故存在  $i_1, i_2, \dots, i_j$ , 使得

$$|\lambda_{i_k} - \mu_k| \leq \|\beta_j \mathbf{q}_{j+1} \mathbf{e}_j^*\| \leq \beta_j. \text{ 证毕.}$$

在 (17) 两边左乘  $Q_j^*$ , 因为  $Q_j$  各列正交于  $\mathbf{q}_{j+1}$ , 因此有

$$Q_j^* A Q_j - Q_j^* Q_j T_j = 0,$$

$$T_j = Q_j^* A Q_j.$$

所以由定理 4.9 可知, 在所有  $j \times j$  矩阵  $M_j$  中, 使  $R = A Q_j - Q_j M_j$  范数  $\|R\|$  最小的是  $T_j$ .

定理 4.12 Ritz 向量  $\mathbf{z}_i = Q_j \mathbf{s}_i$ , 成立如下关系式

$$A \mathbf{z}_i - \mu_i \mathbf{z}_i = \beta_j s_{ji} \mathbf{q}_{j+1},$$

这里  $s_{ji}$  是  $\mathbf{s}_i$  的第  $j$  个分量. 又  $\mathbf{z}_i$  的 Rayleigh 商为  $\mu_i$ , 即

$$\mu_i = (A \mathbf{z}_i, \mathbf{z}_i) / (\mathbf{z}_i, \mathbf{z}_i).$$

证明 在 (17) 两边右乘向量  $\mathbf{s}_i$ ,

$$A Q_j \mathbf{s}_i - Q_j T_j \mathbf{s}_i = \beta_j s_{ji} \mathbf{q}_{j+1},$$

$$A \mathbf{z}_i - \mu_i Q_j \mathbf{s}_i = \beta_j s_{ji} \mathbf{q}_{j+1},$$

$$A \mathbf{z}_i - \mu_i \mathbf{z}_i = \beta_j s_{ji} \mathbf{q}_{j+1},$$



又

$$\begin{aligned}(Az_i, z_i) &= (AQ_j s_i, Q_j s_i) = (Q_j^* A Q_j s_i, s_i) \\ &= (T_j s_i, s_i) = \mu_i,\end{aligned}$$

又  $(z_i, z_i) = (Q_j s_i, Q_j s_i) = (Q_j^* Q_j s_i, s_i) = (s_i, s_i) = 1$ ,

故  $\mu_i = (Az_i, z_i) / (z_i, z_i) = \rho(z_i, A)$ . 证毕.

推论 若  $\lambda$  是  $A$  的特征值中与  $\mu_i$  最接近的一个, 若

$$d_i = \min_{\lambda_l \neq \lambda} |\lambda_l - \mu_i|, \quad (18)$$

$z_i$  与  $y$  的交角为  $\theta_i$ ,  $y$  是对应  $\lambda$  的特征向量, 其意义如定理 4.4, 则

$$|\lambda - \mu_i| \leq (\beta_j s_{ji})^2 / d_i, \quad (19)$$

$$|\sin \theta_i| \leq |\beta_j s_{ji}| / d_i. \quad (20)$$

(19), (20) 两式给出了  $A$  的 Ritz 值和 Ritz 向量, 对于  $A$  的特征值和特征向量逼近程度的估计. 逼近依赖于量  $|\beta_j s_{ji}|$ , 通常将它记为  $\beta_{ji}$ , 如果  $\beta_{ji}$  很小, 那么逼近就会好, 而从分母上的  $d_i$  可以知道, 如果  $\lambda$  是  $A$  的特征值中孤立的一个,  $d_i$  就会大, 逼近就会好. 反之如果  $\lambda$  是  $A$  的特征值的密集丛中的一个  $d_i$  就会小, 于是  $\beta_{ji}/d_i$  就会大, 估计式的右端就大, 实际计算经验也告诉人们, 此时实际逼近也是比较差的.

现在来分析一下  $\beta_{ji}$ , 从

$$r_i = Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1},$$

可以知道  $r_i$  是向量  $Aq_i$  在子空间  $\{q_1, q_2, \dots, q_i\}$  中投影 (最佳逼近) 的偏差, 因为  $\{q_1, q_2, \dots, q_i\}$  就是  $\{q_1, Aq_1, \dots, A^{i-1}q_1\}$ , 因此如果  $q_1$  对于矩阵  $A$  的最小化零多项式是  $m$  次多项式, 则  $q_1, Aq_1, \dots, A^{m-1}q_1$  是线性无关的, 而  $q_1, Aq_1, \dots, A^{m-1}q_1, A^m q_1$  是线性相关的, 于是对于这样的  $q_1$ ,

$$r_i = Aq_i - \alpha_i q_i - \beta_{i-1} q_{i-1} \neq 0, \quad l < m,$$

$$r_m = Aq_m - \alpha_m q_m - \beta_{m-1} q_{m-1} = 0,$$

即  $\beta_l > 0$  当  $1 \leq l < m$ , 而  $\beta_m = 0$ . 此时  $\{q_1, q_2, \dots, q_m\}$  是不变子空间, 对应的 Ritz 值和 Ritz 向量, 恰为  $A$  的特征值和特征向量.

当  $n$  很大时,  $m$  可能也很大, 对于比较大的  $j$ ,  $j < m$ ,  $\beta_j$  是否会很小? 从实际计算的经验来看, 除了少数情况外没有迹象显示  $\beta_j$  会十分小. 而  $|s_{ji}|$  当  $j$  足够大时却常常是十分小的, 尤其是对于  $T_j$  的两端特征值, 所对应的特征向量的  $|s_{ji}|$ , 因此  $\beta_{ji}$  常常可以达到很小, 但这一点理论上还不能给出证明.

下面来讨论  $T_j$  的特征值与  $T_k$  的特征值之间的关系,  $j < k$ . 有

**定理 4.13** 设  $\mu_i^{(j)}$  是  $T_j$  的第  $i$  个特征值, 当  $j < k$  时, 可以找到  $T_k$  的一个特征值  $\mu^{(k)}$ , 使得

$$|\mu_i^{(j)} - \mu^{(k)}| \leq \beta_{ji}.$$

**证明** 设  $s_i$  是  $T_j$  对应  $\mu_i^{(j)}$  的单位特征向量. 造一个  $k$  维列向量  $s'_i = \begin{pmatrix} s_i \\ 0 \end{pmatrix}$ , 也有  $\|s'_i\| = 1$ , 我们来估计

$$\begin{aligned} & \|T_k s'_i - \mu_i^{(j)} s'_i\|, \\ T_k s'_i - \mu_i^{(j)} s'_i &= \begin{pmatrix} T_j & 0 \\ 0 & \beta_j & \\ & \beta_j & \\ & & T_{j+1, k} \end{pmatrix} \begin{pmatrix} s_i \\ 0 \end{pmatrix} - \mu_i^{(j)} \begin{pmatrix} s_i \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} T_j s_i - \mu_i^{(j)} s_i \\ \beta_j s_{ji} \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ \beta_j s_{ji} \\ 0 \end{pmatrix}, \end{aligned}$$

右端是一个只有第  $j+1$  个分量为  $\beta_j s_{ji}$ , 而其余分量都为 0 的向量. 因此

$$\|T_k s'_i - \mu_i^{(j)} s'_i\| = \beta_{ji},$$

利用定理 4.3, 就知存在  $T_k$  的特征值  $\mu^{(k)}$  使得

$$|\mu_i^{(j)} - \mu^{(k)}| \leq \beta_{\mu}. \text{ 证毕.}$$

这个定理告诉我们, 如果计算到某步  $j$ ,  $T_j$  的某个特征值  $\mu_i^{(j)}$ , 它的  $\beta_{\mu}$  很小, 那么  $\mu_i^{(j)}$  是  $A$  的特征值  $\lambda$  很好的近似值, 往后  $T_k$  的特征值中, 一定也有一个很好地逼近  $\lambda$ . 同时也告诉我们, 如果  $T_j$  的特征值  $\mu^{(j)}$  与  $T_k$  的所有特征值都差别很大, 那么  $\mu^{(j)}$  不会是  $A$  的特征值的近似值. 因此这个定理, 对实际计算有指导意义.

我们以  $L^j(A, q_1) = (Q_j, T_j)$  表示对于对称矩阵  $A$ , 单位向量  $q_1$ , 通过 Lanczos 算法计算到  $j$  步, 产生出矩阵  $Q_j = [q_1, q_2, \dots, q_j]$  和  $T_j$ .

**定理 4.14** 如果  $L^j(A, q_1) = (Q_j, T_j)$ , 则

1.  $L^j(\alpha A, q_1) = (Q_j, \alpha T_j)$  对所有正数  $\alpha$ ;
2.  $L^j(A - \alpha I, q_1) = (Q_j, T_j - \alpha I)$  对所有实数  $\alpha$ ;
3.  $L^j(P^*AP, P^*q_1) = (P^*Q_j, T_j)$  对所有  $n \times n$  正交阵  $P$ .

**证明** 1. 记  $B = \alpha A$ , 由  $B, q_1$  通过 Lanczos 算法, 得到  $\tilde{q}_2, \tilde{q}_3, \dots, \tilde{q}_j, \tilde{\alpha}_1, \tilde{\beta}_1, \dots, \tilde{\alpha}_{j-1}, \tilde{\beta}_{j-1}, \tilde{\alpha}_j$ , 于是有

$$\tilde{\beta}_1 \tilde{q}_2 = Bq_1 - \tilde{\alpha}_1 q_1,$$

其中  $\tilde{\alpha}_1 = (Bq_1, q_1) = \alpha(Aq_1, q_1) = \alpha\alpha_1$ .

$$\tilde{\beta}_1 = \|Bq_1 - \tilde{\alpha}_1 q_1\| = \alpha \|Aq_1 - \alpha_1 q_1\| = \alpha\beta_1,$$

于是  $\tilde{\beta}_1 \tilde{q}_2 = Bq_1 - \tilde{\alpha}_1 q_1$ ,

即为  $\alpha\beta_1 \tilde{q}_2 = \alpha(Aq_1 - \alpha_1 q_1) = \alpha\beta_1 q_2$ ,

由此可知  $\tilde{q}_2 = q_2$ . 同理可以证明  $\tilde{q}_3, \dots, \tilde{q}_j$  依次为  $q_3, \dots, q_j$ ; 而  $\tilde{\alpha}_2, \tilde{\beta}_2, \dots, \tilde{\alpha}_{j-1}, \tilde{\beta}_{j-1}, \tilde{\alpha}_j$  依次为  $\alpha\alpha_2, \alpha\beta_2, \dots, \alpha\alpha_{j-1}, \alpha\beta_{j-1}, \alpha\alpha_j$ , 即得  $L^j(\alpha A, q_1) = (Q_j, \alpha T_j)$ .

2 和 3 的证明类同, 不再重复. 证毕.

### § 3 Kaniel-Paige-Saad 理论

在上节分析中知道:  $T_j$  的特征值可以作为  $A$  的某个特征值的近似值, 但不知道是  $A$  的第几个特征值的近似值. S. Kaniel 在 1966 年发表的文章[25]中, 给出了  $\mu_i - \lambda_i$  的估计, 不过证明中有些错误, O. Paige 在 1971 年的学位论文[13]中, 重新研究了这个问题, 给出了  $\mu_i - \lambda_i$  的估计, 并且也给出了  $\sin \phi_i$  的估计, 这里  $\phi_i$  表示 Ritz 向量  $z_i$  与  $A$  的特征向量  $y_i$  的交角. Y. Saad 在 1980 年发表的文章[26]中, 对于  $\mu_i - \lambda_i$  给出比 Paige 的更简单的估计. 所以这些结果, 我们统称为 Kaniel-Paige-Saad 理论.

本节大部分定理的证明思想取自 B. N. Parlett 的书[10].

首先来进一步研究  $A$  的从  $T_j$  得到的 Ritz 向量,  $z_1, z_2, \dots, z_j$  的性质.

#### 1. 正交性

$$(z_i, z_l) = \delta_{il}. \quad (21)$$

因为  $z_i = Q_j s_i$ ,  $z_l = Q_j s_l$ , 故

$$(z_i, z_l) = (Q_j s_i, Q_j s_l) = (Q_j^* Q_j s_i, s_l) = (s_i, s_l) = \delta_{il}.$$

$$2. \quad \{q_1, q_2, \dots, q_j\} = \{z_1, z_2, \dots, z_j\}. \quad (22)$$

记  $[s_1, s_2, \dots, s_j] = S$ , 有  $\det(S) \neq 0$ , 从

$$Q_j S = [z_1, z_2, \dots, z_j],$$

可知(22)成立.

#### 3. $A$ -正交性

$$(Az_l, z_k) = \begin{cases} \mu_l, & l=k, \\ 0, & l \neq k. \end{cases}$$

这是因为

$$\begin{aligned}(Az_i, z_k) &= (Q_j^* A Q_j s_i, s_k) \\ &= (T_j s_i, s_k) = \mu_i (s_i, s_k) = \mu_i \delta_{ik}.\end{aligned}$$

4. 记  $P_j$  表示  $j$  次多项式全体,  $m \in P_j$ , 则

$$(m(A)q_1, z_i) = 0$$

的充要条件是  $(\lambda - \mu_i) \mid m(\lambda)$ .

证明 先证充分性. 假设  $(\lambda - \mu_i) \mid m(\lambda)$ , 有

$$m(\lambda) = (\lambda - \mu_i)g(\lambda), \quad g(\lambda) \in P_{j-1},$$

$$(m(A)q_1, z_i) = (g(A)q_1, (A - \mu_i)z_i),$$

$$g(A)q_1 \in \{q_1, Aq_1, \dots, A^{j-1}q_1\} = \{q_1, q_2, \dots, q_j\},$$

故

$$g(A)q_1 = Q_j h,$$

$h$  是一个  $j$  维向量. 故

$$\begin{aligned}(m(A)q_1, z_i) &= (Q_j h, (A - \mu_i)Q_j s_i) \\ &= (h, (T_j - \mu_i)s_i) = 0.\end{aligned}$$

再证必要性. 若  $(m(A)q_1, z_i) = 0$ ,

$$m(\lambda) = (\lambda - \mu_i)g(\lambda) + d,$$

于是

$$\begin{aligned}(m(A)q_1, z_i) &= ((A - \mu_i I)g(A) + dI)q_1, z_i) \\ &= ((A - \mu_i)g(A)q_1, z_i) + d(q_1, z_i) \\ &= d(q_1, Q_j s_i) = d(e_1, s_i) = ds_{1i},\end{aligned}$$

但由第 2 章的定理 2.9 知  $S_i$  的第一个分量  $s_{1i} \neq 0$ , 故必须  $d=0$ .

类似于 Ritz 向量的性质 4, 对于特征向量  $y_i$  也有如下性质.

5. 若  $(q_1, y_i) \neq 0$ ,  $m(\lambda) \in P_j$ , 则

$$(m(A)q_1, y_i) = 0$$

的充分必要条件是

$$\lambda - \lambda_i \mid m(\lambda).$$

证明 充分性. 若  $\lambda - \lambda_l | m(\lambda)$ , 即  $m(\lambda) = (\lambda - \lambda_l)g(\lambda)$ ,  $g(\lambda) \in P_{j-1}$ , 由此

$$\begin{aligned}(m(A)\mathbf{q}_1, \mathbf{y}_l) &= ((A - \lambda_l I)g(A)\mathbf{q}_1, \mathbf{y}_l) \\ &= (g(A)\mathbf{q}_1, (A - \lambda_l I)\mathbf{y}_l) = 0.\end{aligned}$$

必要性. 若  $(m(A)\mathbf{q}_1, \mathbf{y}_l) = 0$ , 如果  $m(\lambda) = (\lambda - \lambda_l)g(\lambda) + d$ , 则

$$\begin{aligned}((A - \lambda_l I)g(A)\mathbf{q}_1 + d\mathbf{q}_1, \mathbf{y}_l) \\ = ((A - \lambda_l I)g(A)\mathbf{q}_1, \mathbf{y}_l) + d(\mathbf{q}_1, \mathbf{y}_l) = d(\mathbf{q}_1, \mathbf{y}_l) = 0,\end{aligned}$$

但  $(\mathbf{q}_1, \mathbf{y}_l) \neq 0$ , 故必须  $d = 0$ , 即  $(\lambda - \lambda_l) | m(\lambda)$ .

**定理 4.15** 设  $\mu_l$  是  $T_j$  的从小到大排列的第  $l$  个特征值,  $\lambda_l$  是  $A$  的从小到大排列的第  $l$  个特征值, 则

$$\mu_l - \lambda_l \geq 0, \quad \lambda_{n-l+1} - \mu_{j-l+1} \geq 0, \quad l = 1, 2, \dots, j.$$

证明 将  $Q_j = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j]$  扩充成一个正交矩阵  $Q = (Q_j, \tilde{Q}_{n-j})$ ,

$$\begin{aligned}Q^*AQ &= \begin{pmatrix} Q_j^*AQ_j & Q_j^*A\tilde{Q}_{n-j} \\ \tilde{Q}_{n-j}^*AQ_j & \tilde{Q}_{n-j}^*A\tilde{Q}_{n-j} \end{pmatrix} \\ &= \begin{pmatrix} T_j & Q_j^*A\tilde{Q}_{n-j} \\ \tilde{Q}_{n-j}^*AQ_j & \tilde{Q}_{n-j}^*A\tilde{Q}_{n-j} \end{pmatrix},\end{aligned}$$

记  $B = Q^*AQ$ , 它的特征值与  $A$  的相同, 利用极大极小原理知

$$\lambda_l = \min_{V_l} \max_{\mathbf{x} \in V_l} (B\mathbf{x}, \mathbf{x}) / (\mathbf{x}, \mathbf{x}),$$

特别取一个  $l$  维子空间  $\tilde{V}_l$ , 它是由最后  $n-j$  个分量为 0 的向量构成的, 因为  $l \leq j$ , 因此这样的子空间是存在的, 当  $l < j$  时, 这样的子空间可以有很多, 但它们的全体总是一般  $l$  维子空间的全体构成的集合中的子集. 因此

$$\begin{aligned}\lambda_l &\leq \min_{\tilde{V}_l} \max_{\mathbf{x} \in \tilde{V}_l} (B\mathbf{x}, \mathbf{x}) / (\mathbf{x}, \mathbf{x}) \\ &= \min_{\tilde{V}_l} \max_{\mathbf{y} \in \tilde{V}_l} (T_j\mathbf{y}, \mathbf{y}) / (\mathbf{y}, \mathbf{y}) = \mu_l.\end{aligned}$$

同样可证

$$\lambda_{n-i+1} - \mu_{j-i+1} \geq 0. \text{ 证毕.}$$

下面介绍 Kaniel-Paige 理论.

**定理 4.16** 设  $\phi_i = \angle(\mathbf{y}_i, \mathbf{z}_i)$ ,  $i=1, 2, \dots, j$ , 则

$$\sin^2 \phi_i \leq \left[ \mu_i - \lambda_i + \sum_{k=1}^{i-1} (\lambda_{i+1} - \lambda_k) \sin^2 \phi_k \right] / (\lambda_{i+1} - \lambda_i). \quad (23)$$

证明 
$$\mathbf{z}_i = \sum_{k=1}^n r_{ik} \mathbf{y}_k, \quad \sum_{k=1}^n r_{ik}^2 = 1,$$

$$(\mathbf{y}_i, \mathbf{z}_i) = r_{ii} = \cos \phi_i.$$

$$\sin^2 \phi_i = \sum_{\substack{k=1 \\ k \neq i}}^n r_{ik}^2,$$

$$\mu_i = \mathbf{z}_i^* A \mathbf{z}_i = \sum_{k=1}^n r_{ik}^2 \lambda_k,$$

$$\lambda_i = \sum_{k=1}^n r_{ik}^2 \lambda_i,$$

所以

$$\mu_i - \lambda_i = \sum_{\substack{k=1 \\ k \neq i}}^n r_{ik}^2 (\lambda_k - \lambda_i),$$

或

$$\begin{aligned} \mu_i - \lambda_i + \sum_{k=1}^{i-1} r_{ik}^2 (\lambda_i - \lambda_k) &= \sum_{k=i+1}^n r_{ik}^2 (\lambda_k - \lambda_i) \\ &\geq (\lambda_{i+1} - \lambda_i) \sum_{k=i+1}^n r_{ik}^2, \end{aligned}$$

不等式两边都加上项  $(\lambda_{i+1} - \lambda_i) \sum_{k=1}^{i-1} r_{ik}^2$ , 即得

$$\sin^2 \phi_i \leq \left[ \mu_i - \lambda_i + \sum_{k=1}^{i-1} r_{ik}^2 (\lambda_{i+1} - \lambda_k) \right] / (\lambda_{i+1} - \lambda_i), \quad (24)$$

另一方面, 将特征向量  $\mathbf{y}_k$  在 Ritz 向量  $\mathbf{z}_k$  和跟  $\mathbf{z}_k$  的正交方向上作分解,

$$\mathbf{y}_k = \mathbf{z}_k \cos \phi_k + \mathbf{h}_k \sin \phi_k, \quad k=1, 2, \dots, j,$$

$\mathbf{h}_k$  是与  $\mathbf{z}_k$  正交方向上的单位向量, 当  $i \neq k$  时有

$$r_{ik} = (\mathbf{z}_i, \mathbf{y}_k) = (\mathbf{z}_i, \mathbf{y}_k - \mathbf{z}_k \cos \phi_k) = (\mathbf{z}_i, \mathbf{h}_k \sin \phi_k),$$

从而知道  $r_{ik}^2 \leq \sin^2 \phi_k$ ,

将此代入(24)式, 即得

$$\sin^2 \phi_i \leq \left[ \mu_i - \lambda_i + \sum_{k=1}^{i-1} (\lambda_{i+1} - \lambda_k) \sin^2 \phi_k \right] / (\lambda_{i+1} - \lambda_i).$$

证毕.

**定理 4.17** 对任意自然数  $l \leq j$ , 任意非零向量  $\mathbf{s} \in \{\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1\}$ , 如果满足

$$\mathbf{s}^* \mathbf{y}_k = 0, \quad k=1, 2, \dots, l-1,$$

$$\begin{aligned} \text{则} \quad \mu_l &\leq \rho(\mathbf{s}) + \sum_{k=1}^{l-1} (\lambda_n - \mu_k) \sin^2 \phi_k \\ &\leq \rho(\mathbf{s}) + \sum_{k=1}^{l-1} (\lambda_n - \lambda_k) \sin^2 \phi_k, \end{aligned}$$

$$\text{其中} \quad \rho(\mathbf{s}) = (A\mathbf{s}, \mathbf{s}) / (\mathbf{s}, \mathbf{s}).$$

**证明** 结论中第二个不等式成立是明显的, 只要利用定理 4.15 即得.

将  $\mathbf{s}$  分解成

$$\mathbf{s} = \mathbf{t} + \sum_{k=1}^{l-1} r_k \mathbf{z}_k,$$

其中  $\mathbf{t}$  是与  $\mathbf{z}_k (k=1, 2, \dots, l-1)$  正交的向量.

$$(\mathbf{s}, \mathbf{s}) = (\mathbf{t}, \mathbf{t}) + \sum_{k=1}^{l-1} r_k^2,$$

$$r_k = (\mathbf{s}, \mathbf{z}_k) = (\mathbf{s}, \mathbf{z}_k - \mathbf{y}_k \cos \phi_k), \quad k=1, 2, \dots, l-1,$$

$$\text{所以} \quad r_k^2 \leq \|\mathbf{s}\|^2 \sin^2 \phi_k,$$

$$\begin{aligned} (A\mathbf{s}, \mathbf{s}) &= (A\mathbf{t}, \mathbf{t}) + \left( \sum_{k=1}^{l-1} r_k A\mathbf{z}_k, \mathbf{t} \right) \\ &\quad + \left( A\mathbf{t}, \sum_{k=1}^{l-1} r_k \mathbf{z}_k \right) + \left( \sum_{k=1}^{l-1} r_k A\mathbf{z}_k, \sum_{s=1}^{l-1} r_s \mathbf{z}_s \right), \end{aligned}$$



因为  $t \in \{z_1, z_2, \dots, z_l\}$ , 故利用  $z_l$  的性质 3 知

$$(At, z_k) = (Az_k, t) = 0, \quad k=1, 2, \dots, l-1,$$

因此

$$\begin{aligned} \rho(s) &= \left[ (At, t) + \sum_{k=1}^{l-1} r_k^2 (Az_k, z_k) \right] / (s, s) \\ &= \left[ (At, t) + \sum_{k=1}^{l-1} \mu_k r_k^2 \right] / (s, s), \end{aligned}$$

两边减去  $\lambda_n$ , 利用  $(t, t) + \sum_{k=1}^{l-1} r_k^2 = (s, s)$ , 得

$$\begin{aligned} \rho(s) - \lambda_n &= \left[ ((A - \lambda_n I)t, t) + \sum_{k=1}^{l-1} (\mu_k - \lambda_n) r_k^2 \right] / (s, s), \\ ((A - \lambda_n I)t, t) &\leq 0, \quad (s, s) \geq (t, t), \end{aligned}$$

故

$$\rho(s) - \lambda_n \geq \frac{((A - \lambda_n I)t, t)}{(t, t)} + \left( \sum_{k=1}^{l-1} (\mu_k - \lambda_n) r_k^2 \right) / (s, s),$$

因为  $t$  正交于  $z_1, z_2, \dots, z_{l-1}$ , 故

$$\frac{((A - \lambda_n I)t, t)}{(t, t)} \geq \mu_l - \lambda_n,$$

由此得到

$$\mu_l - \lambda_n \leq \rho(s) - \lambda_n + \left( \sum_{k=1}^{l-1} (\lambda_n - \mu_k) r_k^2 \right) / (s, s),$$

再将  $r_k^2 \leq \|s\|^2 \sin^2 \phi_k$  代入, 即得

$$\mu_l \leq \rho(s) + \sum_{k=1}^{l-1} (\lambda_n - \mu_k) \sin^2 \phi_k. \quad \text{证毕.}$$

推论 设非零向量  $s \in \{q_1, Aq_1, \dots, A^{j-1}q_1\}$ ,  $(s, y_k) = 0$ ,  $k=1, 2, \dots, l-1$ ,

$$\rho(s, A - \lambda_l I) = ((A - \lambda_l I)s, s) / (s, s),$$

则

$$\mu_l - \lambda_l \leq \rho(s, A - \lambda_l I) + \sum_{k=1}^{l-1} (\lambda_n - \lambda_k) \sin^2 \phi_k. \quad (25)$$

证明 只要在定理中用  $A - \lambda_l I$  代替  $A$ , 即得.

从定理 4.16 的 (23) 可知如果  $\mu_k - \lambda_k (k=1, 2, \dots, l-1)$  是很小的量, 那么  $\sin^2 \phi_k (k=1, 2, \dots, l-1)$  也是很小的量, 当然这里假定  $\lambda_{k+1} - \lambda_k$  都不是十分小, 在此再利用 (25) 式, 如果  $\rho(s, A - \lambda_l I)$  也是很小的话, 那么  $\mu_l - \lambda_l$  也是小量. 因此要对  $\rho(s, A - \lambda_l I)$  来进行分析.

**定理 4.18** 设  $q_1 = \sum_{i=1}^n \alpha_i y_i$ , 对  $l \leq j$ ,  $\alpha_l \neq 0$ , 记

$$h = \sum_{i=l+1}^n \alpha_i y_i / \left( \sum_{i=l+1}^n \alpha_i^2 \right)^{\frac{1}{2}},$$

则对任意  $m(\lambda) \in P_{l-1}$ , 有

$$\begin{aligned} & \rho(m(A)q_1, A - \lambda_l I) \\ & \leq (\lambda_n - \lambda_l) \left[ \frac{\left( \sum_{i=l+1}^n \alpha_i^2 \right)^{\frac{1}{2}}}{\alpha_l} \frac{\|m(A)h\|}{m(\lambda_l)} \right]^2. \end{aligned} \quad (26)$$

证明

$$\begin{aligned} q_1 &= \sum_{i=1}^l \alpha_i y_i + h \left( \sum_{i=l+1}^n \alpha_i^2 \right)^{\frac{1}{2}}, \\ m(A)q_1 &= \sum_{i=1}^l \alpha_i m(\lambda_i) y_i + m(A)h \left( \sum_{i=l+1}^n \alpha_i^2 \right)^{\frac{1}{2}}, \\ ((A - \lambda_l I)m(A)q_1, m(A)q_1) \\ &= \sum_{i=1}^l \alpha_i^2 m(\lambda_i)^2 (\lambda_i - \lambda_l) \\ &\quad + ((A - \lambda_l I)m(A)h, m(A)h) \left( \sum_{i=l+1}^n \alpha_i^2 \right), \end{aligned}$$

因为  $\lambda_i - \lambda_l \leq 0$ , 因此

$$\begin{aligned} \rho(m(A)q_1, (A - \lambda_l I)) &\leq ((A - \lambda_l I)m(A)h, m(A)h) \\ &\quad \times \left( \sum_{i=l+1}^n \alpha_i^2 \right) / (m(A)q_1, m(A)q_1), \end{aligned}$$

但

$$\begin{aligned} (m(A)q_1, m(A)q_1) &\geq (\alpha_l m(\lambda_l))^2, \\ ((A - \lambda_l I)m(A)h, m(A)h) &\leq (\lambda_n - \lambda_l) \|m(A)h\|^2, \end{aligned}$$

合并这些结果, 即得(26). 证毕.

下面给出  $\mu_l - \lambda_l$  的估计.

**定理 4.19** 对于  $l \leq j$

$$\begin{aligned} \mu_l - \lambda_l &\leq (\lambda_n - \lambda_l) \frac{\sum_{k=l+1}^n \alpha_k^2}{\alpha_l^2} \left( \frac{\prod_{k=1}^{l-1} \frac{\lambda_k - \lambda_n}{\lambda_k - \lambda_l}}{T_{j-l}(1+2r)} \right)^2 \\ &\quad + \sum_{k=1}^{l-1} (\lambda_n - \lambda_k) \sin^2 \phi_k, \end{aligned} \quad (27)$$

这里  $T_k(x)$  是  $k$  阶 Чебышев 多项式,  $r = \frac{\lambda_l - \lambda_{l+1}}{\lambda_{l+1} - \lambda_n}$ .

**证明** 定理 4.17 的推论告诉我们, 只要向量  $s$ , 满足  
1.  $s \in \{q_1, Aq_1, \dots, A^{j-1}q_1\}$ , 2.  $s^*y_k = 0, k=1, 2, \dots, l-1$ ,  
就有

$$\mu_l - \lambda_l \leq \rho(s, A - \lambda_l I) + \sum_{k=1}^{l-1} (\lambda_n - \lambda_k) \sin^2 \phi_k,$$

使用定理 4.18, 取  $s = m(A)q_1$ , 自然满足条件 1, 为了使它满足条件 2, 利用 Ritz 向量的性质 5, 知必须取

$$m(\lambda) = (\lambda - \lambda_1)(\lambda - \lambda_2) \cdots (\lambda - \lambda_{l-1})g(\lambda),$$

其中  $g(\lambda)$  是  $j-l$  次多项式. 对于这样的多项式, 定理 4.17 中的

$$\|m(A)h\| \leq (\lambda_n - \lambda_1)(\lambda_n - \lambda_2) \cdots (\lambda_n - \lambda_{l-1}) \|g(A)h\|,$$

$$m(\lambda_l) = (\lambda_l - \lambda_1)(\lambda_l - \lambda_2) \cdots (\lambda_l - \lambda_{l-1})g(\lambda_l),$$

$h \in \{y_{l+1}, y_{l+2}, \dots, y_n\}$ , 因此若  $h = \sum_{i=l+1}^n \beta_i y_i$ , 有

$$g(A)h = \sum_{i=l+1}^n \beta_i g(\lambda_i) y_i,$$

$$\|g(A)\mathbf{h}\|^2 = \sum_{i=l+1}^n \beta_i^2 (g(\lambda_i))^2,$$

我们取一个多项式  $g(\lambda)$ , 使得  $g(\lambda_i)^2 (i=l+1, \dots, n)$  是小的, 而  $g(\lambda_l)$  是很大的, 那么就能得到好的估计. 取

$$g(\lambda) = T_{j-l} \left( \frac{2\lambda - \lambda_{l+1} - \lambda_n}{\lambda_n - \lambda_{l+1}} \right)$$

就是这样的一个, 此时

$$g(\lambda_i)^2 \leq 1 \quad (i=l+1, l+2, \dots, n),$$

故  $\|g(A)\mathbf{h}\| \leq 1$ , 而  $g(\lambda_l)^2 = T_{j-l}(1+2r)^2$ ,  $r = \frac{\lambda_{l+1} - \lambda_l}{\lambda_n - \lambda_{l+1}}$  随  $r$  增大而增大很快. 对于这样  $g(\lambda)$ , 有

$$\begin{aligned} \mu_l - \lambda_l &\leq (\lambda_n - \lambda_l) \left( \sum_{k=l+1}^n \alpha_k^2 / \alpha_l^2 \right) \\ &\times \left( \prod_{k=1}^{l-1} \frac{\lambda_k - \lambda_n}{\lambda_k - \lambda_l} / T_{j-l}(1+2r) \right)^2 \\ &+ \sum_{k=1}^{l-1} (\lambda_n - \lambda_k) \sin^2 \phi_k. \text{ 证毕.} \end{aligned}$$

利用(23)式和(27)式, 有时可以得到  $\mu_l - \lambda_l$  的估计. 例如: 已知  $\lambda_1=0$ ,  $\lambda_2=0.01$ ,  $\lambda_3=0.04$ ,  $\lambda_4, \lambda_5, \dots, \lambda_n$  都在  $[0.1, 1]$  中, 而

$$\alpha_l = \mathbf{q}_1^* \cdot \mathbf{y}_l = 0.01, \quad l=1, 2, 3,$$

取  $j=53$ .

对于  $l=1$ ,

$$\sum_{k=l+1}^n \alpha_k^2 = 1 - (\alpha_1)^2 = 1 - 0.0001 = 0.9999,$$

$$r = \frac{\lambda_2 - \lambda_1}{\lambda_n - \lambda_2} \leq \frac{0.01}{0.99} = \frac{1}{99},$$

$$T_{j-l}(1+2r) = T_{52} \left( 1 + \frac{2}{99} \right) = 1.73 \times 10^4,$$

$$\mu_1 - \lambda_1 \leq \frac{0.9999}{0.0001} \left[ \frac{1}{T_{52}(1+2r)} \right]^2 \leq 3.34 \times 10^{-5},$$

$$\sin^2 \phi_1 \leq \frac{\mu_1 - \lambda_1}{\lambda_2 - \lambda_1} \leq 3.34 \times 10^{-3}.$$

对于  $l=2$ ,

$$\sum_{k=l+1}^n \alpha_k^2 = 1 - \alpha_1^2 - \alpha_2^2 = 0.9998,$$

$$r = \frac{\lambda_3 - \lambda_2}{\lambda_n - \lambda_3} \leq \frac{0.03}{0.96} = 0.03125,$$

$$T_{51}(1+2 \times 0.03125) = 3.39 \times 10^7,$$

$$\mu_2 - \lambda_2 \leq \frac{0.9998}{0.0001} \times \left[ \frac{1}{0.01} / T_{51}(1+0.0625) \right]^2$$

$$+ \sin^2 \phi_1 \leq 3.35 \times 10^{-3},$$

$$\begin{aligned} \sin^2 \phi_2 &\leq \frac{\mu_2 - \lambda_2}{\lambda_3 - \lambda_2} + \sin^2 \phi_1 \frac{\lambda_3 - \lambda_1}{\lambda_3 - \lambda_2} \\ &= \frac{3.35 \times 10^{-3}}{0.03} + 3.34 \times 10^{-3} \frac{0.04}{0.03} \\ &= 1.16 \times 10^{-1}. \end{aligned}$$

依此类推, 可以获得  $\mu_l - \lambda_l$  和  $\sin^2 \phi_l$  的估计. 当然这种估计, 离开实用还是很远, 因为  $\alpha_l$  不易知道,  $\lambda_l$  也是不知道的. 但是从理论上可以看到, 当  $j \rightarrow \infty$  时,  $\mu_l - \lambda_l$  和  $\sin^2 \phi_l$  都会趋于 0, 因此是有理论价值的. 同时也告诉我们: Lanczos 方法, 首先求出来的是端部的特征值, 因为  $\mu_l - \lambda_l$  和  $\mu_{l+1} - \lambda_{l+1}$ , 分别依赖于  $T_{j-l}(1+2r)$  和  $T_{j-l-1}(1+2r)$ , 而前者常比后者大.

上面 Kaniel-Paige 的估计, 还可以稍许改进, 有人考虑如下: 在估计式 (27) 中  $T_{j-l}(1+2r)$ , 随  $j$  增大而增大的速度, 依赖于  $r = \frac{\lambda_l - \lambda_{l+1}}{\lambda_{l+1} - \lambda_n}$ , 如果矩阵  $A$  的特征值分布如图 4 所示,

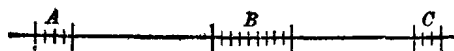


图 4

分成三部分, 左端部分  $A$ , 中间部分  $B$ , 右端部分  $C$ ,  $A$  和  $C$  与中间部分  $B$  却分隔得较远.

区间  $A$  有特征值  $\lambda_1, \lambda_2, \dots, \lambda_{r_1}$ ,

区间  $B$  有特征值  $\lambda_{r_1+1}, \lambda_{r_1+2}, \dots, \lambda_{r_2}$ ,

区间  $C$  有特征值  $\lambda_{r_2+1}, \dots, \lambda_n$ ,

$$\text{此时取 } m(\lambda) = \frac{\prod_{k=1}^{r_1} (\lambda - \lambda_k)}{(\lambda - \lambda_l)} \prod_{k=r_2+1}^n (\lambda - \lambda_k) g(\lambda),$$

$$g(\lambda) = T_{j-(n-r_2)-r_1} \left( \frac{2\lambda - \lambda_{r_1+1} - \lambda_{r_2}}{\lambda_{r_2} - \lambda_{r_1+1}} \right),$$

$$g(\lambda_l)^2 = T_{j-(n-r_2)-r_1} \left( \frac{2\lambda_l - \lambda_{r_1+1} - \lambda_{r_2}}{\lambda_{r_2} - \lambda_{r_1+1}} \right)^2$$

$$= T_{j-(n-r_2)-r_1} (1 + 2r_{r_1 r_2})^2,$$

$$r_{r_1 r_2} = \frac{\lambda_{r_1+1} - \lambda_l}{\lambda_{r_2} - \lambda_{r_1+1}},$$

这样可使  $r_{r_1 r_2}$  比较大, 从而会获得更好的估计.

现在开始介绍 Saad 的结果; Paige 为了获得  $\mu_l - \lambda_l$  的估计考虑的是  $s \in \{q_1, Aq_1, \dots, A^{j-1}q_1\}$ , 并且  $s$  正交于特征向量  $y_1, y_2, \dots, y_{l-1}$ . Saad 考虑的是  $s \in \{q_1, Aq_1, \dots, A^{j-1}q_1\}$ ,  $s$  正交于 Ritz 向量  $z_1, z_2, \dots, z_{l-1}$ .

**定理 4.20** 对任意  $s \in \{q_1, Aq_1, \dots, A^{j-1}q_1\}$ , 如果  $s$  正交于 Ritz 向量  $z_1, z_2, \dots, z_{l-1}$ , 则

$$\mu_l - \lambda_l \leq \rho(s, A - \lambda_l I).$$

**证明** 只要证明

$$\mu_l \leq \rho(s, A),$$

因为  $\mathbf{s} \in \{\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1\}$ , 故  $\mathbf{s} \in \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_j\}$ , 因此

$$\mathbf{s} = \sum_{k=1}^j \beta_k \mathbf{z}_k,$$

但  $\mathbf{s}^* \mathbf{z}_k = 0 (k=1, 2, \dots, l-1)$ , 故

$$\mathbf{s} = \sum_{k=l}^j \beta_k \mathbf{z}_k = Q_j \sum_{k=l}^j \beta_k \mathbf{s}_k,$$

$$\rho(\mathbf{s}, A) = (A\mathbf{s}, \mathbf{s}) / (\mathbf{s}, \mathbf{s})$$

$$= \left( Q_j^* A Q_j \sum_{k=l}^j \beta_k \mathbf{s}_k, \sum_{k=l}^j \beta_k \mathbf{s}_k \right) / \left( \sum_{k=l}^j \beta_k \mathbf{s}_k, \sum_{k=l}^j \beta_k \mathbf{s}_k \right) \\ = (T_j \mathbf{h}, \mathbf{h}) / (\mathbf{h}, \mathbf{h}) = \rho(\mathbf{h}, T_j),$$

这里  $\mathbf{h} = \sum_{k=l}^j \beta_k \mathbf{s}_k$ , 是正交于  $T_j$  的特征向量  $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_{l-1}$  的

$j$  维向量, 因此  $\rho(\mathbf{h}, T_j) \geq \mu_l$ , 即

$$\mu_l \leq \rho(\mathbf{s}, A). \text{ 证毕.}$$

从定理 4.20 可以看出 Saad 结果会简单一些.

**定理 4.21** 对任意自然数  $l \leq j \leq n$ ,

$$\mathbf{q}_1 = \sum_{k=1}^n \alpha_k \mathbf{y}_k, \quad \alpha_l \neq 0,$$

如果  $\mu_{l-1} \leq \lambda_{l+1}$ , 则

$$\mu_l - \lambda_l \leq (\lambda_n - \lambda_l) \left( \sum_{k=l+1}^n \alpha_k^2 / \alpha_l^2 \right) \\ \times \left( \prod_{k=1}^{l-1} \frac{\mu_k - \lambda_n}{\mu_k - \lambda_l} / T_{j-l}(1+2r) \right)^2, \quad (28)$$

这里

$$r = \frac{\lambda_l - \lambda_{l+1}}{\lambda_{l+1} - \lambda_n}.$$

**证明** 利用定理 4.18, 对任意  $m(\lambda) \in P_{j-1}$ , 有

$$\rho(m(A)\mathbf{q}_1, A - \lambda_l I) \leq (\lambda_n - \lambda_l) \left( \sum_{k=l+1}^n \alpha_k^2 / \alpha_l^2 \right) \frac{\|m(A)\mathbf{h}\|^2}{m(\alpha_l)^2},$$

$\mathbf{s} = m(A)\mathbf{q}_1 \in \{\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1\}$ , 由 Ritz 向量的性质 4 知, 要  $\mathbf{s}^* \mathbf{z}_k = 0 (k=1, 2, \dots, l-1)$ , 可取

$$m(\lambda) = (\lambda - \mu_1)(\lambda - \mu_2) \cdots (\lambda - \mu_{l-1})g(\lambda),$$

其中  $g(\lambda)$  是任意  $j-l$  次多项式, 对于这样的  $m(\lambda)$ , 利用定理 4.20, 有

$$\mu_l - \lambda_l \leq \rho(s, A - \lambda_l I) \leq (\lambda_n - \lambda_l) \left( \sum_{k=l+1}^n \alpha_k^2 / \alpha_l^2 \right) \frac{\|m(A)h\|^2}{m(\alpha_l)^2},$$

但

$$\begin{aligned} \|m(A)h\| &\leq \max_{l+1 \leq i \leq n} |(\lambda_l - \mu_1)(\lambda_l - \mu_2) \cdots \\ &\quad \cdot (\lambda_l - \mu_{l-1})| \max_{l+1 \leq i \leq n} |g(\lambda_i)| \\ &= (\lambda_n - \mu_1)(\lambda_n - \mu_2) \cdots (\lambda_n - \mu_{l-1}) \max_{l+1 \leq i \leq n} |g(\lambda_i)|, \end{aligned}$$

特别取 
$$g(\lambda) = T_{j-l} \left( \frac{2\lambda - \lambda_{l+1} - \lambda_n}{\lambda_n - \lambda_{l+1}} \right),$$

即得(28). 证毕.

对于 Saad 的结果, 也可以取

$$\begin{aligned} m(\lambda) &= \frac{\prod_{k=1}^{r_1} (\lambda - \lambda_k)}{(\lambda - \lambda_l)} \\ &\quad \times \prod_{k=r_2+1}^n (\lambda - \lambda_k) T_{j-(n-r_2)-r_1} \left( \frac{2\lambda - \lambda_{r_1+1} - \lambda_{r_2}}{\lambda_{r_2} - \lambda_{r_1+1}} \right) \end{aligned}$$

来改善  $\mu_l - \lambda_l$  的估计.

## § 4 在有限位精度运算下的 Lanczos 算法

在有限位精度运算下的 Lanczos 算法实施时, 会产生舍入误差, 它的影响破坏了  $q_l$  与  $q_k$  之间的正交性, 于是当  $j$  适当大时,  $I - Q_j^* Q_j$  不再是零矩阵, 甚至  $\|I - Q_j^* Q_j\|$  不是一个小量. 与前述无限位精度运算下的结果相比, 带来了复杂性.



例如, 对于

$$\beta_{l-1} = \|r_{l-1}\|,$$

$$r_l = Aq_l - \alpha_l q_l - \beta_{l-1} q_{l-1},$$

在无限精度下:

$$\beta_{l-1} = (Aq_l, q_{l-1}).$$

但在有限精度下  $\|r_{l-1}\|$ , 与  $(Aq_l, q_{l-1})$  就不一定相同.

在无限精度下

$$\alpha_l = (Aq_l, q_l) = (Aq_l - \beta_{l-1} q_{l-1}, q_l),$$

但在有限精度下  $(Aq_l, q_l)$  与  $(Aq_l - \beta_{l-1} q_{l-1}, q_l)$  也不一定相同.

从这两个例子说明, 对于每一步计算中,  $\alpha_l$  与  $\beta_{l-1}$  按不同公式计算是有差别的. 为了下面讨论方便, 我们规定, Lanczos 算法如下:

1. 取向量  $q_1$ ,  $\|q_1\| = 1$ , 造  $u_1 = Aq_1$ ,  $1 \rightarrow l$ ;
2.  $\alpha_l = q_l^* u_l$ ;
3.  $r_l = u_l - \alpha_l q_l$ ;
4.  $\beta_l = \|r_l\|$ ;
5. 若  $\beta_l$  充分小, 转入计算  $T_l$  的特征值, 特征向量. 否则  $q_{l+1} = r_l / \beta_l$ ;
6.  $u_{l+1} = Aq_{l+1} - \beta_l q_l$ ,  $l+1 \rightarrow l$  转 2.

这样算法, 可使矩阵  $T_l$ , 保持对称, 使得

$$(r_l, q_l) = 0$$

在工作精度下成立. 但是  $(r_l, q_{l-1})$  不一定是零.

在这样的算法下, 关系式

$$AQ_j - Q_j T_j = \beta_j q_{j+1} e_j^*$$

也不一定成立了. 但可以知道成立

$$AQ_j - Q_j T_j = \beta_j q_{j+1} e_j^* + F_j, \quad (29)$$

这里  $F_j$  是由舍入误差引起的  $n \times j$  矩阵。

从实际计算的实验知

$$\|F_j\| \leq \varepsilon \|A\|,$$

这里  $\varepsilon$  是计算机的工作精度。在 C. Paige 1976 的文章 [27] 中, 指出

$$\|F_j\|_F \leq (7 + \alpha) \sqrt{j} \varepsilon \|A\|$$

一定成立。这里  $\alpha$  是跟计算  $Ax$  有关的误差。如果  $A$  的每行最多有  $m$  个非零元, 将  $A$  的元素都取绝对值后的矩阵记为  $|A|$ ,  $\||A|\|/\|A\| = \nu$ , 则

$$\alpha \leq m\nu.$$

关于这一结果, 我们不再介绍了。我们的兴趣, 仍然是估计  $T_j$  的特征值与  $A$  的特征值之间的差, 也即要回答  $T_j$  的特征值, 是否可以作为  $A$  的特征值的近似值。

因为现在  $Q_j$  的列不再正交了, 因此我们可以利用 W. Kahan 定理 4.10, 得到

$$|\lambda_{i_k} - \mu_k| \leq \sqrt{2} \|R\| / \sigma_1,$$

这里  $R = \beta_j q_{j+1} e_j^* + F_j$ , 来估计  $\mu_k$  作为  $A$  的特征值的近似值。可是计算经验表明, 当  $q_1, q_2, \dots, q_j$ ,  $j$  增大时, 不但失去正交性, 而且很快就线性相关了。因此在  $j$  不太大时, 可以用上述估计式。当  $j$  较大时,  $\sigma_1(Q_j) = 0$ , 就不能使用这个估计式了。为此我们转而研究 (7) 型的估计式。

假如  $T_j$  的特征值  $\mu_k$  所对应的特征向量为  $s_k$ ,  $\|s_k\| = 1$ , 从 (29) 式有

$$AQ_j s_k - \mu_k Q_j s_k = \beta_j s_{jk} q_{j+1} + F_j s_k,$$

仍记  $z_k = Q_j s_k$  为 Ritz 向量, 于是有

$$Az_k - \mu_k z_k = \beta_j s_{jk} q_{j+1} + F_j s_k,$$

于是利用 (7) 有  $A$  的特征值  $\lambda$ ,

$$|\lambda - \mu_k| \leq \|r\| / \|z_k\|,$$

这里  $r = \beta_j s_{jk} q_{j+1} + F_j s_k$ , 有

$$\|r\| \leq \beta_{jk} + \|F_j\|.$$

因为  $\|F_j\|$  是一个小量, 因此只要  $\beta_{jk}$  充分小, 可使  $\|r\|$  是一个小量. 留下的问题是  $\|z_k\|$  是否会很小.

尽管  $\|s_k\| = 1$ , 但  $Q_j$  的列不正交, 甚至会线性相关, 因此  $Q_j s_k$  的范数可能会很小.

这样使用 Lanczos 算法, 面临二个问题:

1.  $\beta_{jk}$  是否会很小, 什么时候它才很小.
2.  $\|z_k\|$  是否能保持在 1 附近.

1971 年 O. Paige 的博士论文 [13] 回答了这两个问题.

他的结论是

1.  $z_k$  与  $q_{j+1}$  失去正交性  $\iff \beta_{jk}$  很小;
2. 当  $\min_{i \neq k} |\mu_i - \mu_k|$  不很小时,  $\|z_k\|$  不会很小.

下面来介绍这些结果, 记

$$I - Q_j^* Q_j = C_j^* + \Delta_j + C_j,$$

这里  $C_j$  是严格上三角形,  $\Delta_j$  是对角阵, 在我们规定的 Lanczos 算法下,  $(\tilde{q}_l, q_{l+1}) = 0$ , 因此我们有  $C_j$  的第 1 条超对角上的元素都为 0.  $\Delta_j$  的对角元反映  $\|q_l\| = 1$  的误差, 因此一般说来也是很小的.

记  $s_1, s_2, \dots, s_j$  是  $T_j$  正确的单位特征向量, 对应的特征值为  $\mu_1, \mu_2, \dots, \mu_j$ , 因此矩阵

$$S = [s_1, s_2, \dots, s_j]$$

是一个正交阵. 矩阵  $F_j^* Q_j - Q_j^* F_j$  是一个斜对称矩阵, 因此它的对角元为 0, 故可记为

$$F_j^* Q_j - Q_j^* F_j = K - K^*,$$

$K$  是严格上三角阵, 由于  $\|F_j\| \sim \varepsilon \|A\|$ , 因此  $K$  一般也不会太大, 记  $s_i^*(K - K^*)s_k = l_{ik}^{(j)}$ .

$$\Delta_j T_j - T_j \Delta_j$$

也是一个斜对称矩阵, 记为  $N - N^*$ ,  $N$  是严格上三角阵, 因为  $\Delta_j$  的元素是很小的量, 因此  $N$  的元素也是小量.

矩阵  $C_j T_j - T_j C_j$  是一个对角元为 0 的上三角阵. 记矩阵

$$G = S^*(K + N)S = (r_{ik}^{(j)}).$$

**定理 4.22** 设  $q_i^* q_{l+1} = 0$ ,  $l = 1, 2, \dots, j$ , Ritz 向量  $z_i = Q_j s_i$ ,  $i = 1, 2, \dots, j$ , 则

$$1. \quad z_i^* q_{j+1} = r_{ii}^{(j)} / \beta_j s_{jj}, \quad i = 1, 2, \dots, j; \quad (30)$$

2. 当  $i \neq k$  时

$$(\mu_i - \mu_k) z_i^* z_k = r_{ii}^{(j)} \frac{\beta_j s_{jk}}{\beta_j s_{jj}} - r_{kk}^{(j)} \frac{\beta_j s_{ji}}{\beta_j s_{jj}} - l_{ik}^{(j)}. \quad (31)$$

**证明** 由 (29) 式

$$A Q_j - Q_j T_j = \beta_j q_{j+1} e_j^* + F_j,$$

$$\text{得} \quad Q_j^* A Q_j - Q_j^* Q_j T_j = \beta_j Q_j^* q_{j+1} e_j^* + Q_j^* F_j,$$

两边取转置

$$Q_j^* A Q_j - T_j Q_j^* Q_j = \beta_j e_j q_{j+1}^* Q_j + F_j^* Q_j,$$

两式相减得

$$T_j Q_j^* Q_j - Q_j^* Q_j T_j = \beta_j Q_j^* q_{j+1} e_j^* - \beta_j e_j q_{j+1}^* Q_j + Q_j^* F_j - F_j^* Q_j, \quad (32)$$

即

$$\begin{aligned} & (I - Q_j^* Q_j) T_j - T_j (I - Q_j^* Q_j) \\ & = \beta_j Q_j^* q_{j+1} e_j^* - \beta_j e_j q_{j+1}^* Q_j + K^* - K, \end{aligned} \quad (33)$$

等式左边为

$$\begin{aligned}
& (C_j^* + \Delta_j + C_j)T_j - T_j(C_j^* + \Delta_j + C_j) \\
& = C_jT_j - T_jC_j + \Delta_jT_j - T_j\Delta_j + C_j^*T_j - T_jC_j^* \\
& = C_jT_j - T_jC_j + N - N^* + C_j^*T_j - T_jC_j^*,
\end{aligned}$$

于是比较(33)式等式两边的上三角部分, 得到

$$\begin{aligned}
C_jT_j - T_jC_j + N &= \beta_j Q_j^* q_{j+1} e_j^* - K, \\
\beta_j Q_j^* q_{j+1} e_j^* &= C_jT_j - T_jC_j + N + K,
\end{aligned}$$

两边左乘  $s_i^*$ , 右乘  $s_i$  得到

$$\beta_j s_{ji} z_i^* q_{j+1} = \mu_i s_i^* C_j s_i - \mu_i s_i^* C_j s_i + r_{ii}^{(j)} = r_{ii}^{(j)},$$

$$\text{即 } z_i^* q_{j+1} = r_{ii}^{(j)} / \beta_j s_{ji} \quad (j=1, 2, \dots, j),$$

此即第一个结论.

再在(32)两边, 左乘  $s_i^*$ , 右乘  $s_k$

$$\begin{aligned}
(\mu_i - \mu_k) z_i^* z_k &= \beta_j s_{jk} z_i^* q_{j+1} - \beta_j s_{ji} q_{j+1}^* z_k \\
&+ s_i^* (Q^* F_j - F_j^* Q_j) s_k \\
&= \beta_j s_{jk} r_{ii}^{(j)} / \beta_j s_{ji} - \beta_j s_{ji} r_{kk}^{(j)} / \beta_j s_{jk} - l_{ik}^{(j)},
\end{aligned}$$

即为第二个结论. 证毕.

从定理 4.22 知道, 当  $z_i$  与  $q_{j+1}$  正交性好的时候, 即  $z_i^* q_{j+1} \sim 0$  时,  $\beta_{ji}$  不能很小, 反之当  $\beta_{ji}$  很小时  $z_i$ 、 $q_{j+1}$  正交性就差.

从(31)可以知道当  $\mu_i \neq \mu_k$  时, 要  $z_i^* z_k$  正交性好,  $\beta_j s_{jk}$  与  $\beta_j s_{ji}$  的数量级应该差不多, 如果两者数量级差别很大, 就使  $|z_i^* z_k|$  不是很小.

下面的定理, 给出  $\|z_k\|$  的估计.

**定理 4.23** 在前述假定下, 如记

$$\nu_k = \min_{i \neq k} |\mu_i - \mu_k|, \quad \text{若 } \nu_k \neq 0,$$

$$\zeta_k = j(j-1)r/\nu_k,$$

$$r = \max_{s, l} |r_{ss}^{(l)}|,$$

则

$$|1 - \|z_k\|| \leq \|A_j\| + \zeta_k. \quad (34)$$

$$\begin{aligned} \text{证明 } 1 - \|z_k\|^2 &= 1 - z_k^* z_k = 1 - s_k^* Q_j^* Q_j s_k \\ &= s_k^* (I - Q_j^* Q_j) s_k = s_k^* (C_j^* + A_j + C_j) s_k, \end{aligned}$$

因为

$$\|s_k^* C_j s_k\| = \|s_k^* C_j^* s_k\|,$$

故

$$|1 - \|z_k\|^2| \leq \|A_j\| + 2\|s_k^* C_j s_k\|.$$

由此可知: 要获得定理 4.23 的结果(34)式, 主要的任务是估计  $\|s_k^* C_j s_k\|$ , 为此应用定理 4.22 结果给出以下二个引理.

**引理 4.1** 记  $T_i$ ,  $i < j$  的标准正交特征向量组为  $s_1^{(i)}$ ,  $s_2^{(i)}$ ,  $\dots$ ,  $s_i^{(i)}$ , 正交阵  $S_i = [s_1^{(i)}, s_2^{(i)}, \dots, s_i^{(i)}]$ , 又记  $j \times i$  长方形

$$S'_i = \begin{pmatrix} S_i \\ 0 \end{pmatrix}$$

按照定理 4.22 所用的记号, 记向量

$$v_i = \left( \frac{r_{1,1}^{(i)}}{\beta_i s_{i,1}^{(i)}}, \frac{r_{2,2}^{(i)}}{\beta_i s_{i,2}^{(i)}}, \dots, \frac{r_{i,i}^{(i)}}{\beta_i s_{i,i}^{(i)}} \right)^T, \quad i = 1, 2, \dots, j-1,$$

则

$$s_k^* C_j s_k = - \sum_{l=1}^{j-1} (s_k^* S'_l v_l) s_{l+1,k},$$

这里  $s_{l+1,k}$  是  $T_l$  的特征向量  $s_k$  的第  $l+1$  个分量.

**证明**

$$-C_j = \begin{pmatrix} 0 & q_1^* q_2 & q_1^* q_3 & \cdots & q_1^* q_j \\ & 0 & q_2^* q_3 & \cdots & q_2^* q_j \\ & & 0 & \ddots & \vdots \\ & & & \ddots & q_{j-1}^* q_j \\ & & & & 0 \end{pmatrix},$$

所以  $-C_j$  的第  $i+1$  列为

$$\begin{pmatrix} Q_i^* q_{i+1} \\ 0 \end{pmatrix} = \begin{pmatrix} S_i S_i^* Q_i^* q_{i+1} \\ 0 \end{pmatrix},$$

若记  $Z_i = Q_i S_i$ , 它是  $T_i$  的 Ritz 向量组成的矩阵, 这样  $-C_i$  的第  $i+1$  列, 又可写成

$$\begin{pmatrix} S_i Z_i^* q_{i+1} \\ 0 \end{pmatrix},$$

对于  $Z_i^* q_{i+1}$  可以应用定理 4.22, 有

$$Z_i^* q_{i+1} = \left( \frac{r_{1,1}^{(i)}}{\beta_i s_{i,1}^{(i)}}, \frac{r_{2,2}^{(i)}}{\beta_i s_{i,2}^{(i)}}, \dots, \frac{r_{i,i}^{(i)}}{\beta_i s_{i,i}^{(i)}} \right)^T = v_i,$$

于是  $-C_i$  的第  $i+1$  列为

$$\begin{pmatrix} S_i v_i \\ 0 \end{pmatrix} = S_i' v_i,$$

因此  $s_k^* C_i s_k = - \sum_{l=1}^{i-1} (s_k^* S_l' v_l) s_{i+1, k}$ . 证毕.

现在需要估计  $s_k^* S_l'$ .

**引理 4.2** 若记  $T_j$  的特征值为  $\mu_1, \mu_2, \dots, \mu_j, T_l (l < j)$  的特征值为  $\mu_1^{(l)}, \mu_2^{(l)}, \dots, \mu_l^{(l)}$ , 都是按由小到大次序排列. 记

$$\omega_{l,k} = \begin{cases} \sum_{p=1}^l r_{p,p}^{(i)} / (\mu_k - \mu_p^{(i)}), & \text{当 } \mu_k \neq \mu_p^{(i)}, p=1, 2, \dots, l, \\ 0, & \text{当 } \mu_k \text{ 与 } \mu_p^{(i)} (p=1, 2, \dots, l) \\ & \text{中某一个相同时.} \end{cases}$$

则  $s_k^* C_j s_k = - \sum_{l=1}^{j-1} s_{i+1, k}^2 \omega_{l, k}$ . (35)

证明  $S_l'$  中的第  $p$  列为

$$\begin{pmatrix} s_p^{(i)} \\ 0 \end{pmatrix},$$

$s_p^{(i)}$  是  $T_l$  的第  $p$  个特征向量, 对应的特征值为  $\mu_p^{(i)}$ ,

$$T_j \begin{pmatrix} s_p^{(i)} \\ 0 \end{pmatrix} - \mu_p^{(i)} \begin{pmatrix} s_p^{(i)} \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ \beta_i s_{i,p}^{(i)} \\ 0 \end{pmatrix} \text{ 第 } l+1 \text{ 行,}$$

两边左乘  $\mathbf{s}_k^*$  得

$$\mu_k \mathbf{s}_k^* \begin{pmatrix} \mathbf{s}_p^{(l)} \\ 0 \end{pmatrix} - \mu_p^{(l)} \mathbf{s}_k^* \begin{pmatrix} \mathbf{s}_p^{(l)} \\ 0 \end{pmatrix} = \beta_l s_{l,p}^{(l)} s_{l+1,k},$$

或 
$$(\mu_k - \mu_p^{(l)}) \mathbf{s}_k^* \begin{pmatrix} \mathbf{s}_p^{(l)} \\ 0 \end{pmatrix} = \beta_l s_{l,p}^{(l)} s_{l+1,k},$$

若  $\mu_k - \mu_p^{(l)} \neq 0$ , 则

$$\mathbf{s}_k^* \begin{pmatrix} \mathbf{s}_p^{(l)} \\ 0 \end{pmatrix} = \frac{\beta_l s_{l,p}^{(l)} s_{l+1,k}}{(\mu_k - \mu_p^{(l)})},$$

如果对某个  $p$ , 有  $\mu_k - \mu_p^{(l)} = 0$ , 那么必须有

$$\beta_l s_{l,p}^{(l)} s_{l+1,k} = 0,$$

但因为  $l < j$ ,  $\beta_l \neq 0$ ,  $s_{l,p}^{(l)}$  是  $\mathbf{s}_p^{(l)}$  的最后一个分量, 在第 2 章已经指出, 它不为 0, 因此必须  $s_{l+1,k} = 0$ . 这样如果我们记某些自然数的集合

$$m = \{l \mid \{1, 2, \dots, j-1\} \cap \{\mu_k \neq \mu_r^{(l)}, r=1, 2, \dots, l\},$$

则有

$$\begin{aligned} \mathbf{s}_k^* C_j \mathbf{s}_k &= - \sum_{l \in m} (\mathbf{s}_k^* S_l^T \mathbf{v}_l) s_{l+1,k} \\ &= - \sum_{l \in m} (\mathbf{s}_k^* \mathbf{s}_1^{(l)}, \mathbf{s}_k^* \mathbf{s}_2^{(l)}, \dots, \mathbf{s}_k^* \mathbf{s}_l^{(l)}) \mathbf{v}_l s_{l+1,k} \\ &= - \sum_{l \in m} \beta_l s_{l+1,k}^2 \left( \frac{s_{l,1}^{(l)}}{\mu_k - \mu_1^{(l)}}, \frac{s_{l,2}^{(l)}}{\mu_k - \mu_2^{(l)}}, \dots, \frac{s_{l,l}^{(l)}}{\mu_k - \mu_l^{(l)}} \right) \mathbf{v}_l \\ &= - \sum_{l \in m} \beta_l s_{l+1,k}^2 \sum_{p=1}^l \frac{s_{l,p}^{(l)} r_{p,p}^{(l)}}{(\mu_k - \mu_p^{(l)}) \beta_l s_{l,p}^{(l)}} \\ &= - \sum_{l \in m} s_{l+1,k}^2 \sum_{p=1}^l \frac{r_{p,p}^{(l)}}{\mu_k - \mu_p^{(l)}}, \end{aligned}$$

因此若记 
$$\omega_{l,k} = \begin{cases} \sum_{p=1}^l \frac{r_{p,p}^{(l)}}{\mu_k - \mu_p^{(l)}}, & l \in m, \\ 0, & l \notin m, \end{cases}$$

(这就是本引理结论中的量  $\omega_{l,k}$ ) 就有



$$s_k^* C_j s_k = - \sum_{i=1}^{l-1} s_{i+1,k}^2 \omega_{l,k}. \text{ 证毕.}$$

用引理 4.2, 如果我们假定

$$\tilde{\nu}_k = \min_{\substack{\mu_k^{(j)} + \mu_k \\ l=1, 2, \dots, j-1 \\ r=1, 2, \dots, l}} |\mu_k - \mu_r^{(l)}|,$$

那么我们可以很容易获得类似于定理的结果. 但是  $\tilde{\nu}_k$  这个量可能很小很小, 因此这样获得的结果是粗糙的, 为了获得更好的结果, 将  $\tilde{\nu}_k$  换成  $\nu_k$ . 下面要想法把  $\mu_k - \mu_r^{(l)}$  换成  $\mu_k - \mu_{l_i}$ .

现在再来继续证明定理 4.23, 考虑量  $s_{i+1,k}^2$ , 由第 2 章定理 2.9 公式 (15),

$$s_{i+1,k}^2 = \chi_{1,i}(\mu_k) \chi_{l+2,j}(\mu_k) / \chi'_{1,j}(\mu_k),$$

$\chi_{1,i}(\lambda)$  是  $T_i$  的特征多项式, 它的根为  $\mu_1^{(i)}, \mu_2^{(i)}, \dots, \mu_i^{(i)}$ ,  $\chi_{l+2,j}(\lambda)$  有  $j-l-1$  个根. 将  $\chi_{1,i}(\lambda)$  和  $\chi_{l+2,j}(\lambda)$  的全部  $j-1$  个根, 按从小到大次序排列成  $\omega_1, \omega_2, \dots, \omega_{k-1}, \omega_{k+1}, \dots, \omega_j$ ,

因为  $\chi'_{1,j}(\mu_k) = \prod_{\substack{i=1 \\ i \neq k}}^j (\mu_k - \mu_i)$ , 因此

$$s_{i+1,k}^2 = \prod_{i=1}^{k-1} \frac{\mu_k - \omega_i}{\mu_k - \mu_{l_i}} \prod_{i=k+1}^j \frac{\mu_k - \omega_i}{\mu_k - \mu_{l_i}}.$$

因为  $\omega_1, \omega_2, \dots, \omega_{k-1}, \omega_{k+1}, \dots, \omega_j$  是矩阵  $T_j$  划去第  $l+1$  行、第  $l+1$  列后的矩阵  $B_{l+1}$  的特征值, 利用极大极小原理可知

$$\left. \begin{aligned} \mu_s &\leq \omega_s \quad (s=1, 2, \dots, k-1), \\ \mu_k &\leq \omega_{k+1}, \\ \mu_{j-s} &\geq \omega_{j-s} \quad (s=0, 1, \dots, j-k-1), \\ \mu_k &\geq \omega_{k-1}. \end{aligned} \right\} \quad (36)$$

实际上, 如果对于第  $l+1$  个分量为 0 的向量  $x$ , 有

$$\rho(x, T_j) = \rho(x, B_{l+1}),$$

若记  $\tilde{V}_s$  是  $j$  维向量的一个  $s$  维子空间, 这些向量的第  $l+1$  个分量为 0. 于是

$$\begin{aligned}\mu_s &= \min_{V_s} \max_{x \in V_s} \rho(x, T_j) \leq \min_{\tilde{V}_s} \max_{x \in \tilde{V}_s} \rho(x, T_j) \\ &= \min_{\tilde{V}_s} \max_{x \in \tilde{V}_s} \rho(x, B_{l+1}) \\ &= B_{l+1} \text{ 的第 } s \text{ 个特征值} \quad (s \leq j-1).\end{aligned}$$

同样

$$\begin{aligned}\mu_{j-s+1} &= \max_{V_s} \min_{x \in V_s} \rho(x, T_j) \geq \max_{\tilde{V}_s} \min_{x \in \tilde{V}_s} \rho(x, T_j) \\ &= \max_{\tilde{V}_s} \min_{x \in \tilde{V}_s} \rho(x, B_{l+1}) \\ &= B_{l+1} \text{ 的第 } j-s+1 \text{ 个特征值}, \quad s \leq j-1.\end{aligned}$$

由(36)式可知

$$\left| \frac{\mu_k - \omega_i}{\mu_k - \mu_i} \right| \leq 1 \quad (i=1, 2, \dots, k-1, k+1, k+2, \dots, j). \quad (37)$$

另一方面对每个  $r \leq l$ , 在  $\omega_1, \omega_2, \dots, \omega_{k-1}, \omega_{k+1}, \dots, \omega_j$  中必有一个是  $\mu_r^{(l)}$ , 记为第  $i(r)$  个, 即  $\mu_r^{(l)} = \omega_{i(r)}$ , 并且  $i(r) \neq k$ , 于是

$$\begin{aligned}s_{l+1,k}^2 &= \prod_{i=1}^{k-1} \frac{\mu_k - \omega_i}{\mu_k - \mu_i} \prod_{i=k+1}^j \frac{\mu_k - \omega_i}{\mu_k - \mu_i} \\ &= g_{i(r)} \frac{\mu_k - \mu_r^{(l)}}{\mu_k - \mu_{i(r)}} \quad (r=1, 2, \dots, l),\end{aligned}$$

其中

$$g_{i(r)} = s_{l+1,k}^2 / \frac{\mu_k - \omega_{i(r)}}{\mu_k - \mu_{i(r)}},$$

由(37)知  $|g_{i(r)}| \leq 1$ .

$$\begin{aligned}\text{由此 } s_k^* O_j s_k &= - \sum_{l=1}^{j-1} s_{l+1,k}^2 \omega_{l,k} \\ &= - \sum_{l=1}^{j-1} \sum_{p=1}^l g_{i(p)} \frac{\mu_k - \mu_p^{(l)}}{\mu_k - \mu_{i(p)}} \frac{\tau_{p,p}^{(l)}}{\mu_k - \mu_p^{(l)}} \\ &= - \sum_{l=1}^{j-1} \sum_{p=1}^l g_{i(p)} \frac{\tau_{p,p}^{(l)}}{\mu_k - \mu_{i(p)}},\end{aligned}$$

从而

$$|s^* C_j s_k| \leq \sum_{i=1}^{j-1} \sum_{p=1}^i \frac{|r_{p,p}^{(j)}|}{|\mu_k - \mu_{i(p)}|} \leq \frac{r}{\nu_k} \sum_{i=1}^{j-1} \sum_{p=1}^i 1 \\ \leq \frac{j(j-1)}{2} r / \nu_k. \quad \text{证毕.}$$

$r_{s,s}^{(j)}$  是跟  $\varepsilon \|A\|$  相当的小量, 因此  $r$  也是跟  $\varepsilon \|A\|$  相当的小量, 故若  $\nu_k$  不是太小,  $\zeta_k$  一般也是一个小量. 于是从定理 4.22 知  $1 - \|z_k\|$ , 在  $\nu_k$  不是太小的条件下, 也是一个小量.

至此我们介绍了 Lanczos 的理论基础. 在实际计算时 Lanczos 算法过程中, 常常会出现这样的问题: 某个  $T_j$  的二个特征值或更多个特征值, 对应逼近的是  $A$  的同一个特征值. 这一现象的出现, 是由于  $Q_j$  的正交性失掉而产生的. 在计算中如果碰到这种情况, 就要把这些特征值的 Ritz 向量算出来, 看看是否平行或几乎平行. 如果是平行或几乎平行, 那么这些 Ritz 值就对应  $A$  的同一个特征值.

从上面的分析可知,  $\mu_k$  是否接近于  $A$  的特征值, 关键是看量  $\beta_{jk}$  是否足够小. 但  $\beta_{jk} = |\beta_j s_{jk}|$ ,  $\beta_j$  是已知的量,  $s_{jk}$  是  $T_j$  的第  $k$  个特征向量的最后一个分量. 如果能提供一种方法, 直接从  $T_j$  计算得  $s_{jk}$ , 对 Lanczos 算法是很有价值的.

下面我们来介绍一种求  $T_j$  的特征向量的办法, 假定特征值  $\mu_k$  已经知道了.

$$(T_j - \mu_k I)z = 0, \quad (38)$$

$z = (\eta_1, \eta_2, \dots, \eta_j)^T$ , 令  $\eta_1 = 1$ , 理论上可以从 (38) 的第 1 个方程求得  $\eta_2$ , 再从第二个方程求得  $\eta_3$ , 依此类推可以从 (38) 的第  $j-1$  个方程算得第  $j$  个分量  $\eta_j$ . 但是这样的方法在 Wilkinson 的书 [20] 中, 已经分析过, 是不稳定的. 原因是我们使用的是近似特征值  $\tilde{\mu}_k$ , 计算又有舍入误差, 因此  $\eta_{j-1}$  和  $\eta_j$  两个量代入到 (38) 的第  $j$  个方程左边,

甲

四  
四  
一  
七

$$(\alpha_j - \mu_k)\eta_j + \beta_{j-1}\eta_{j-1}$$

不会是 0, 如果它是  $\delta$ , 于是得到方程

$$(T_j - \tilde{\mu}_k I)z = \delta e_j,$$

这样 
$$z = \delta \sum_{l=1}^j \frac{(e_j, s_l)}{\mu_l - \mu_k} s_l = \delta \sum_{l=1}^j \frac{s_{jl}}{\mu_l - \mu_k} s_k,$$

如果  $\frac{s_{jk}}{\mu_k - \mu_k}$  的数量级比其余  $\frac{s_{jl}}{\mu_l - \mu_k}$  的数量级要大, 那么可以使  $z$  是  $s_k$  方向好的近似向量. 否则  $z$  就不是好的近似特征向量.

同样从  $\eta_j = 1$ , 利用方程 (38) 朝前推, 求得  $\eta_{j-1}, \eta_{j-2}, \dots, \eta_1$  的方法, 也有同样的问题. 此时得到的方程是

$$(T_j - \tilde{\mu}_k I)z = \delta e_1,$$

B. N. Parlett 和 J. K. Reid 在文 [28] 中, 提出如下的方法: 对  $(T_j - \tilde{\mu}_k I)z = 0$ , 令  $\eta_1 = 1$ , 从第 1 个方程求出  $\eta_2$ , 依次到从第  $l-1$  个方程求出  $\eta_l$ .

再令  $\xi_j = 1$ , 从第  $j$  个方程

$$(\alpha_j - \tilde{\mu}_k)\xi_j + \beta_{j-1}\xi_{j-1} = 0$$

求出  $\xi_{j-1}$ , 依次到从第  $(l+1)$  个方程求得  $\xi_l$ .

再将分别从两端出发求得的两组数啮合起来, 即令

$$\xi_i = \eta_i \frac{\xi_i}{\eta_i} \quad (i=1, 2, \dots, l-1), \quad (39)$$

这样得到的向量  $z = (\xi_1, \xi_2, \dots, \xi_l)^T$ , 代入第  $l$  个方程左边有

$$\beta_{l-1}\xi_{l-1} + (\alpha_l - \tilde{\mu}_k)\xi_l + \beta_l\xi_{l+1} = \delta, \quad (40)$$

于是  $z$  满足  $(T_j - \tilde{\mu}_k I)z = \delta e_l$ .

现在来考虑如何选择指标  $l$ , 使得  $(z, e_l)$  尽可能大, 也即  $z$  的分量中最大的一个对应的次序作为  $l$ . 为此从  $\eta_1 = 1, \eta_2, \dots$  到某  $\eta_m$ , 绝对值都逐个不减少, 但  $|\eta_{m+1}| < |\eta_m|$ , 那么

这个  $m$  就是  $l$  的一个候选者。再从  $\xi_j=1, \xi_{j-1}, \dots$  逐个计算到  $\xi_m$ , 通过啮合公式 (40) 求出  $\mathbf{z}=(\xi_1, \xi_2, \dots, \xi_j)^T$  后, 将它的分量代入 (40), 得到  $\delta$ , 考察  $\delta/\|\mathbf{z}\|$ , 如果这是一个小量, 那么  $\mathbf{s}=\mathbf{z}/\|\mathbf{z}\|$  即为所求的单位特征向量的近似。否则, 再从  $\eta_{m+1}$  开始求第二个极大值  $\eta_{m'}$ , 当然  $|\eta_{m'}|$  要求比  $|\eta_m|$  大, 以  $m'$  作为  $l$  的候选者, 以此类推。

为什么要以判定  $\delta/\|\mathbf{z}\|$  的大小, 作为一个判据? 利用关系式 (29)

$$AQ_j - Q_j T_j = \beta_j \mathbf{q}_{j+1} \mathbf{e}_j^* + F_j,$$

$$\text{和} \quad (T_j - \tilde{\mu}_k I) \mathbf{s} = \delta \mathbf{e}_l / \|\mathbf{z}\|,$$

$$\text{因此} \quad AQ_j \mathbf{s} - Q_j T_j \mathbf{s} = \beta_j \mathbf{q}_{j+1} \mathbf{e}_j^* \mathbf{s} + F_j \mathbf{s},$$

$$\text{有} \quad AQ_j \mathbf{s} - \tilde{\mu}_k Q_j \mathbf{s} = \beta_j \mathbf{q}_{j+1} \mathbf{e}_j^* \mathbf{s} + F_j \mathbf{s} + \frac{\delta}{\|\mathbf{z}\|} Q_j \mathbf{e}_l,$$

记  $\mathbf{y} = Q_j \mathbf{s}$ , 又记  $\mathbf{s}$  最后一个分量为  $\tilde{s}_{jk}$ , 即有

$$A\mathbf{y} - \tilde{\mu}_k \mathbf{y} = \beta_j \tilde{s}_{jk} \mathbf{q}_{j+1} + F_j \mathbf{s} + \frac{\delta}{\|\mathbf{z}\|} \mathbf{q}_l,$$

$\mathbf{y}$  逼近特征向量的程度、 $\tilde{\mu}_k$  逼近特征值的程度, 都依赖于

$$r = \beta_j \tilde{s}_{jk} \mathbf{q}_{j+1} + F_j \mathbf{s} + \frac{\delta}{\|\mathbf{z}\|} \mathbf{q}_l$$

的范数, 而  $r$  的第三项的范数为  $\delta/\|\mathbf{z}\|$ , 因此使用  $\delta/\|\mathbf{z}\|$  作为判据。

利用反迭代法求特征向量, 也是一个稳定的方法, 只是存储量和计算量比上述方法稍许大一点, 参见 [28]。

最后再说一下, 使用 Lanczos 算法一般地  $j$  应该取多大, 根据 Parlett 的经验 [10, p.259], 当  $n=10^4$ , 取  $j=300$ , 可以求出 10 个两端部的特征值和对应的特征向量。要进一步求更多的特征值和特征向量, 可以继续增大  $j$ , 但这样要遇到已求得

的特征值的重复出现的现象。也可以重新取一个初始单位向量  $\tilde{q}_1$ ，令它与已经求得的特征向量正交，从  $\tilde{q}_1$  出发再施行 Lanczos 过程。

适当地修改 Lanczos 过程，防止特征值重复出现的现象，仍然是一个研究课题，参见 [29]、[30]。

## 参 考 文 献

- [1] M. R. Hestenes and E. Stiefel, Methods of conjugate gradients for solving linear systems, NBS J. Res. 49 pp. 409-436, 1952.
- [2] J. K. Reid, On the method of conjugate gradients for the solution of large sparse system of linear equations, Proceedings of the Conference on Large Sparse Sets of Linear Equations (Ed. J. K. Reid), Academic Press, pp. 231-254. 1971.
- [3] G. W. Stewart, The convergence of the method of conjugate gradients at isolated extreme Points of the spectrum, Numer. Math. 24: 85-93 (1975).
- [4] Dianne P. O'Leary, The block conjugate gradient algorithm and related methods, Linear Algebra and Its Applications 29:293-322, 1980.
- [5] 蒋尔雄 轴对称热应力问题的算法和程序, 748 会议论文集, 1974.
- [6] 蒋尔雄 对称正定矩阵  $P$ -条件数的改善, 复旦大学学报, 自然科学版, Vol. 9, No. 1:1-7, 1964.
- [7] J. A. Meigerink and H. A. Van der Vorst, An iterative solution method for linear systems of which the coefficient matrix is a symmetric  $M$ -matrix, Mathematics of Computation 31:148-162, 1977.
- [8] David S. Kershaw, The incomplete Cholesky-conjugate gradient method for the iterative solution of systems of linear equations, Journal of Computational Physics 26, 43-65, 1978.
- [9] 蒋尔雄、高坤敏、吴景琨 线性代数, 人民教育出版社 1978.
- [10] B. N. Parlett, The Symmetric Eigenvalue Problem, Prentice-Hall, Inc. Englewood Cliffs, N. J. 1980.
- [11] C. H. Reinsch, A stable, rational QR algorithm for the computation of the eigenvalues of an Hermitian tridiagonal matrix, Math. of Comp., Vol. 25, No. 115, 1971.
- [12] R. C. Thompson and P. McEntegert, Principal submatrices II, the upper and lower Quadratic inequalities, Linear Algebra and its

Appl., 1:211-243, 1968.

- [13] C. C. Paige, The computation of eigenvalues and eigenvectors of very large sparse matrices, Ph. D. thesis, Univ. of London, 1971.
- [14] Harry Hochstadt, On some inverse problems in matrix theory, Arch. Math. 18:201-207, 1967.
- [15] Ole H. Hald, Inverse eigenvalue problems for Jacobi matrices, Linear Algebra Appl. 14:63-85, 1976.
- [16] C. de Boor and G. H. Golub, The numerically stable reconstruction of a Jacobi matrix from spectral data, Linear Algebra and its Appl. 21:245-260, 1978.
- [17] 曹维藩 关于一种特殊形状线性方程组的求解问题, 计算数学, 1978 年第 1 期.
- [18] C. C. Paige and Saunders, Solution of Sparse indefinite systems of linear equations, SIAM J. Numer. Anal. 12:617-629, 1975.
- [19] B. N. Parlett, A new look at the Lanczos algorithm for solving symmetric systems of linear equations, Linear Algebra and its Applications. 29:323-346, 1980.
- [20] J. H. Wilkinson, The Algebraic Eigenvalue Problem, Clarendon Press. Oxford, 1965.
- [21] J. H. Wilkinson, Global convergence of tridiagonal QR algorithm with origin shifts, J. Linear Algebra and Its Applications. Vol. 1:409-420, 1968.
- [22] C. L. Lawson and R. J. Hanson, Solving Least Squares Problems, Prentice-Hall, INC., 1974.
- [23] C. Lanczos, An iteration method for the solution of the eigenvalue problem of linear differential and integral operators, J. Res. Nat. Bur. Standards, Sect. B 45:225-280, 1950.
- [24] H. R. Schwarz, Tridiagonalization of a symmetric band matrix, Numer. Math. 12:231-241, 1968.
- [25] S. Kaniel, Estimates for some computational techniques in linear algebra, Math. Comp. 20:369-378, 1966.
- [26] Y. Saad, Error bounds on the interior Rayleigh-Ritz approximations from Krylov subspaces, Soc. Ind. Appl. Math. J. Num. Anal. 17, 1980.
- [27] C. Paige, Error analysis of Lanczos algorithm for tridiagonalizing a



- symmetric matrix, J. Inst. Math. Appl. 18:341-349, 1976.
- [28] B. N. Parlett and J. K. Reid, Tracking the progress of the Lanczos algorithm for large symmetric eigenproblems, IMA Jour Num. Analysis, Vol. 1:135-155, 1981.
- [29] Parlett, B. N. and Scott, D. S. The Lanczos algorithm with selective orthogonalization, Math. Comp. 33:217-238, 1979.
- [30] Horst D. Simon, The Lanczos algorithm for solving symmetric linear systems, Ph. D. Thesis, University of California, Berkeley, 1982.
- [31] 蒋尔雄 对称三对角矩阵带位移的 QL 方法的收敛率, 高等学校计算数学学报, Vol. 6, No. 3, 1984.
- [32] Jiang Erriong and Zhang Zhenyue, A new shift of the QL algorithm for irreducible symmetric tridiagonal matrices, Linear Algebra and Its Applications, 将发表.

[General Information]

书名=对称矩阵计算

作者=蒋尔雄

页数=160

SS号=10184133

DX号=

出版日期=1984年11月第1版

出版社=上海科学技术出版社

封面

书名

版权

前言

目录

## 第1章 共轭斜量法

### § 1 斜量法

### § 2 多步斜量法

### § 3 共轭斜量法

### § 4 不完全分解、预处理共轭斜量法

## 第2章 对称三对角矩阵

### § 1 Jacobi 矩阵

### § 2 对称三对角矩阵的唯一归化定理

### § 3 对称三对角矩阵的极值性质

### § 4 Thompson-McEntegget-Paige公式和特征值反问题

### § 5 解对称线性代数方程组的Lanczos算法

## 第3章 解特征值问题的QL方法

### § 1 QL方法的一般性质

### § 2 用于对称三对角矩阵时的QL方法的性质

### § 3 带Rayleigh商位移的QL方法

### § 4 带Wilkinson位移的QL方法

## 第4章 解特征值问题的Lanczos算法

### § 1 近似不变子空间

### § 2 Lanczos算法

### § 3 Kaniel-Paige-Saad理论

### § 4 在有限位精度运算下的Lanczos算法

参考文献